# Deliverable

## D5.4 Tools for detecting single nucleotide poly-morphisms and analyses within and between hosts

**Version: 1.0**

**Due: Month 36**

**Completed: Month 59**
**Author: Ron A.M. Fouchier (Erasmus MC)**

## Contents

# Deliverable description

An important added value of NGS/WGS methods is the potential to detect SNPs associated with specific pathogen traits, including drug resistance and virulence, not only in consensus pathogen genome sequences, but also as (minor) variants among heterogeneous mixtures of sequences (e.g. quasispecies). Finding such minority variants can advance the detection of relevant changes by days or even weeks. Reducing the signal to noise ratio through pre-analytical steps (WP2) as well as downstream quality controlled bio-informatic analyses are crucial for reliable application of NGS to identification of informative polymorphisms. While current-day practice can lead to quantifying minor SNP variants to ~1% reliably (e.g. Linster et al. Cell 2014), for some applications there is a clear need to develop technologies beyond this threshold (Russell et al., Science 2014). This task will focus on optimizing quantitative output on SNPs by developing appropriate laboratory methods and analysis tools.

# Deliverable execution

With the emergence of Highly Pathogenic Avian Influenza (HPAI) viruses of the H5N8 subtype in the EU around the start of the COMPARE project, it was rapidly decided to (re-)focus several COMPARE tasks to these outbreaks. The unique opportunity to apply NGS research and applications to the imminent thread posed by HPAI H5 viruses for wild birds, poultry and perhaps even human health was considered particularly relevant for the SNP detection tasks. The ability to describe SNPs correctly and reliably was considered crucial for communication between EU laboratories about e.g. the risks associated with some of these SNPs (e.g. virulence determinants, host range determinants, transmission determinants, drugs resistance markers; see Deliverable 5.5) and source attribution. Beyond simple SNPs, there was strong interest also in minor variant analyses. What if the COMPARE phylogeography projects on the global migration of HPAI H5 viruses (see Global Consortium for H5N8 and Related Influenza Viruses. Science. 2016 Oct 14;354(6309):213-217) could have increased resolution based on minor variant analyses? What if deep sequencing and minor variant analyses could improve the relevance or sensitivity of phenotype predictions (see Deliverable 5.5) from (deep) genotyping data?

Three COMPARE partners (APHA, FLI, Erasmus MC) performed NGS on three closely related H5N8 viruses, that were each analyzed using different NGS strategies (Illumina with and without amplification, 454), sequencing instruments (Illumina, 454) and data processing pipelines (three in house versions, partially using commercial tools as well). The goal was to compare the consensus virus genome sequences as well as minor variants therein as determined in the three centers and identify sources of error. In order to determine the comparability of consensus sequences and minority (sub-consensus) single nucleotide variant identification, the biological samples, the sequence data from the three sequencing platforms and the *.bam quality-trimmed alignment files of raw data of the three influenza A/H5N8 viruses were shared among the partners using the EMBL-EBI datahubs.

# Results

To test the applicability of real-time sequence data sharing within the COMPARE network, all raw sequence data used in this study were uploaded to and shared via a datahub in the ENA environment. Using this hub, sharing between institutions was greatly facilitated and immediate access to the data prior to the public release was realized to enable joint evaluation and comparison. All data files have been made publicly available via the ENA (https://www.ebi.ac.uk/ena).

The analyses of consensus virus genome sequences revealed 100% agreement between platforms. Reliable consensus sequences were generated independently of the sequencing platform and data processing pipeline used, although the well-known artefactual InDels in homopolymer regions using the Roche 454 genome sequencer required manual editing. These known problems were not followed up further because this platform was discontinued by the manufacturer anyway. We conclude that consensus sequences used for the detailed characterization of influenza virus strains in outbreak situations can be called reliably with NGS approaches.

In contrast to the reproducible generation of consensus virus genome sequences, we concluded that minority variants were not identified reproducibly. Observed differences were mainly attributed to the alignment processes in the different data processing pipelines and sequencing depth of the sequencing platforms. There was limited reproducibility of minor variant identification data, even for relative high frequency mSNVs. The reproducibility was best (30%) for high frequency (≥10%) variants, and least (9.4% to 31.1%) for the low frequency (≥1%) variants.

We conclude that minority variant analyses will need a different level of careful standardization and awareness about the possible limitations, as shown in this study. Future NGS research projects should address these issues.

## Output

The details for this deliverable will hopefully be published. The draft manuscript is provided with this deliverable. Unfortunately, the manuscript has been in a review process for well over a year, which explains the delays for this deliverable report. In part, this appears to be due to the undesirable outcome of the minor variant analysis for many scientists in the NGS community, possibly including the editor(s) and reviewers in the review process.

# Annex 1 Draft article:

Comparison of sequencing methods and data processing pipelines for whole genome sequencing and minority single nucleotide variant (mSNV) analysis during an influenza A/H5N8 outbreak

# Comparison of sequencing methods and data processing pipelines for whole genome sequencing and minority single nucleotide variant (mSNV) analysis during an influenza A/H5N8 outbreak

Marjolein J. Poen[1], Anne Pohlmann[2], Clara Amid[3], Theo M. Bestebroer[1], Sharon M. Brooks[4], Ian H. Brown[4], Helen Everett[4], Claudia M.E. Schapendonk[1], Rachel D. Scheuer[1], Saskia L. Smits[1], Martin Beer[2], Ron A.M. Fouchier[1], Richard J. Ellis[4*]

1. Erasmus MC, Department of Viroscience, Rotterdam, the Netherlands

2. Institute of Diagnostic Virology, Friedrich-Loeffler-Institute, Insel Riems, Germany

3. European Molecular Biology Laboratory (EMBL), European Bioinformatics Institute (EBI), Wellcome Genome Campus, Hinxton, Cambridge, United Kingdom

4. Animal and Plant Health Agency (APHA) - Weybridge, Addlestone, Surrey, United Kingdom

* corresponding author: richard.ellis@apha.gov.uk

# Abstract

As high-throughput sequencing technologies are becoming more widely adopted for analysing pathogens in disease outbreaks there needs to be assurance that the different sequencing technologies and approaches to data analysis will yield reliable and comparable results. Conversely, understanding where agreement cannot be achieved provides insight into the limitations of these approaches and also allows efforts to be focused on areas of the process that need improvement. This manuscript describes the next-generation sequencing of three closely related viruses, each analysed using different sequencing strategies, sequencing instruments and data processing pipelines. In order to determine the comparability of consensus sequences and minority (sub-consensus) single nucleotide variant (mSNV) identification, the biological samples, the sequence data from 3 sequencing platforms and the *.bam quality-trimmed alignment files of raw data of 3 influenza A/H5N8 viruses were shared. This analysis demonstrated that variation in the final result could be attributed to all stages in the process, but the most critical were the well-known homopolymer errors introduced by 454 sequencing, and the alignment processes in the different data processing pipelines which affected the consistency of mSNV detection. However, homopolymer errors aside, there was generally a good agreement between consensus sequences that were obtained for all combinations of sequencing platforms and data processing pipelines. Nevertheless, minority variant analysis will need a different level of careful standardization and awareness about the possible limitations, as shown in this study.

# Introduction

Over the past decade, high-throughput sequencing technologies have evolved, providing faster, cheaper, and less laborious alternatives to obtain (whole genome) DNA and RNA sequences compared to traditional Sanger sequencing [1, 2]. The use of next-generation sequencing (NGS) technologies is continuously expanding and has revolutionized the field of genomics and molecular biology.

In many fields of infectious disease research, nucleotide changes in DNA or RNA sequences are used to monitor genetic adaptions indicative of evolution, the emergence of drug resistance, immune evasion or as a tool in epidemiological tracing [3]. In clinical settings, sequencing information is used to improve diagnostics and prognosis. NGS technologies play an increasingly important role in these processes as clinically or epidemiologically important nucleotide changes can be present in the minority of DNA or RNA sequences only, which might be missed with more traditional (consensus) sequencing methods which determine the most abundant sequence variants in a population. Nucleotide variants that are present in only a minority of the sequenced virus population are referred to as minority Single Nucleotide Variants (mSNVs). These variants, initially occurring due to replication errors, can become fixed in the population when they have some sort of evolutionary advantage, for instance, mutations related to drug resistance. Furthermore, mSNVs can be also used for high-resolution molecular epidemiology, which becomes more and more important for outbreak assessment [4, 5]. Traditional Sanger sequencing for instance has been described to detect minority variants provided they are present in at least 10% of the analysed DNA or RNA strands within a sample [6, 7]. Hence, the use of traditional sequencing methods is usually restricted to obtaining consensus sequences or to determine heterozygosity in diploid organisms. In contrast, NGS technologies are able to detect low frequency mSNVs in

3

68    sequence fragments or even whole genomes. Typically, NGS sensitivity for minority

69    sequence variant identification is restricted to a level of variation of 0.1–1%, mainly due to

70    sequencing related background errors [8-10], but sensitivity can be increased using

71    sophisticated approaches like circle sequencing [11] or improved bioinformatic analysis

72    workflows [10]. The reliability of mSNV analysis using NGS methods is influenced by many

73    factors, like the quantity and quality of the input sample, the laboratory procedures, the type

74    of sequencing platform and the software and settings used to analyse the raw sequence data.


75    Due to the technical improvements, NGS technologies have become more important as

76    diagnostic tools to characterize pathogens in outbreak situations. However, the increasing use

77    of these technologies to address new and important (outbreak related) research and

78    surveillance questions emphasizes the need to determine the reproducibility of, and the

79    important technical considerations affecting, outcomes obtained by different laboratories

80    following different protocols. Given this, comparative studies focusing on different platforms

81    and data analysis methods are essential to cross-validate different methodologies and

82    determine the reliability of newly obtained data. In addition, there is a growing need (as

83    exemplified by the recent Ebola and Zika virus outbreaks) to share also comprehensive

84    sequencing data as quickly as possible to help with source attribution and developing control

85    strategies. However, the underlying technologies and methods used for NGS are still diverse

86    and there is a strong demand for harmonization of laboratory procedures and approaches for a

87    reliable and optimized analysis of the data.


88    This study is part of the European Union's HORIZON 2020 project "COMPARE"

89    (http://www.compare-europe.eu/), aiming to improve the analytical tools for emerging

90    zoonotic pathogens and its underpinning research. Here, the comparability of NGS output

91    data obtained from different sequence approaches were evaluated and demonstrated suitable

92    sharing strategies for comprehensive NGS data sets. In November 2014, a newly emerging

strain of highly pathogenic avian influenza (HPAI) virus was detected in several European

countries [12, 13]. In the United Kingdom [14], Germany [15], and The Netherlands [16-18]

this subtype was detected in commercial poultry farms within a few days of one another. In

each of those countries, NGS was used to generate whole-genome sequences rapidly after

detection, but as the laboratories in each country were working independently, different

approaches were used for both sequencing and data analysis, and the data were shared as part

of a wider study to determine the likely source of the outbreak [19]. It is important to

determine whether the different analytical approaches have any impact on the outcome.

Therefore, the aim of this study was to determine how comparable consensus and minority

variant results were between laboratories performing their standard analyses, and whether

discrepancies could be attributed to the sequence platform (SP), the data processing platform

(DPP) or a combination of both. With the lack of a ground truth/gold standard, all datasets

obtained were compared amongst each other. The hypothesis we test in this study is that

outputs from NGS analysis of viruses will be comparable irrespective of laboratory,

sequencing platform and data analysis platform.

Therefore, virus isolates obtained in each of the three countries (United Kingdom, Germany

and the Netherlands) were shared between these three partners and subsequently sequenced

and analysed in each of the three laboratories according to local procedures. In addition, the

use of a specially designed data sharing platform, a COMPARE "Data Hub" at EMBL-EBI,

Hinxton UK, was evaluated. This study presents genome coverage data, consensus

sequences, the analysis of the comparability of mSNV identifications of the different SPs,

and DPPs.

Our hypothesis was confirmed at the consensus sequence level, since consensus sequences

could be reproduced independent of the combination of SP and DPP used. However, the

identification of minority variants appeared to be poorly reproducible, primarily due to the

118     well-known errors in 454 sequencing, and due to differences induced by the alignment

119     processes in the different DPPs. The interpretation of minority variant analysis thus needs a

120     different level of careful standardization and awareness about the possible limitations as

121     shown in this study.

122

# Materials and Methods

## Experimental design

125     Three avian influenza A virus isolates that were obtained from three different avian species

126     during the 2014/15 outbreak of HPAI H5N8 virus in Europe were shared among three

127     institutions in the United Kingdom (Animal Plant and Health Agency [APHA]), Germany

128     (Friedrich-Loeffler-Institut [FLI]) and the Netherlands (Erasmus Medical Center [EMC]),

129     later referred to as anonymized institutions I, II and III (Figure 1). All three institutions

130     sequenced all three virus isolates according to their own standard procedures. Adaptors used

131     in the sequencing processes were trimmed off before the raw sequence data files were shared.

132     The sequence data files (*.fastq files), alignment files (*.bam files), sample metadata and

133     experimental metadata were shared between the three laboratories and analysed in their own

134     DPPs yielding sequence datasets for each virus (Table 1). This approach enabled to separate

135     the biological features of the viruses from variation introduced by technical methodology.

136     Data sharing was facilitated via a "Data Hub" provided by the EMBL-EBI's European

137     Nucleotide Archive (ENA) in the framework of the COMPARE collaborative project; all data

138     were stored and subsequently published in ENA [20] (https://www.ebi.ac.uk/ena, for the

139     accession numbers, see Table 1). ENA is an open repository for sequence and related data

140     and a member of the International Nucleotide Sequence Database Collaboration (INSDC;

141 http://www.insdc.org/) [21]. A full description of the COMPARE Data Hub system is

142 provided in a preprint version of Amid et al. [22]. First, consensus sequences derived from a

143 preliminary analysis were compared and one overarching consensus sequence was

144 determined for each gene segment for each virus. This custom-made consensus was used by

145 all three institutions as the reference genome for undertaking mSNV analysis. The resulting

146 nine mSNV reports (originating from three whole-genome raw data sequences times three

147 DPPs) were combined for all three viruses in one spreadsheet file per virus to check the

148 reproducibility of mSNV identification when using different combinations of SP and DPP.

149 The experimental design is summarized in figure 1.

150

151

152 **Table 1. Sample characteristics and accession details**

| | UKDD | | | DETU | | | NLCH | | |
|---|---|---|---|---|---|---|---|---|---|
| **Virus strain** | A/duck/England/36254/2014 | | | A/turkey/Germany/AR2485-L01478/2014 | | | A/chicken/Netherlands/EMC-3/2014 | | |
| **Isolation source** | Pooled intestines | | | Lung tissue | | | Lung tissue | | |
| **Host Scientific Name** | Anas platyrhynchos | | | Meleagris gallopavo | | | Gallus gallus domesticus | | |
| **Host Common Name** | Domestic duck | | | Turkey | | | Chicken | | |
| **Collection Date** | 14 November 2014 | | | 04 November 2014 | | | 23 November 2014 | | |
| **Collection Country** | United Kingdom | | | Germany | | | Netherlands | | |
| **Collection Region** | East Yorkshire | | | Mecklenburg-Western Pomerania | | | Ter Aar | | |
| **Influenza Test Method** | MP gene RRT-PCR, H5 RRT-PCR | | | MP gene RRT-PCR, H5 RRT-PCR | | | MP gene RRT-PCR, H5 RRT-PCR | | |
| **Culture Status Sample** | Egg passage 1 | | | Egg passage 1 | | | MDCK passage 2 | | |
| | **Institution I** | **Institution II** | **Institution III** | **Institution I** | **Institution II** | **Institution III** | **Institution I** | **Institution II** | **Institution III** |
| **Study Accession*** | PRJEB9846 | PRJEB12582 | PRJEB9687 | PRJEB9846 | PRJEB12582 | PRJEB9687 | PRJEB9846 | PRJEB12582 | PRJEB9687 |

7

| | ERR 9728 05 | ERR 1293 054 | ERR 9267 12 | ERR 9267 13 | ERR 1354 020 | ERR 1293 053 | ERR 9267 14 | ERR 9267 15 | ERR 1354 021 | ERR 1293 055 | ERR 9267 17 | ERR 9267 18 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Run Accession*** | | | | | | | | | | | | |
| **DPP1 *.bam file run accession*** | ERR 3093 746 | ERR 3093 752 | ERR903375 6 | | ERR 3093 744 | ERR 3093 753 | ERR309375 7 | | ERR 3093 745 | ERR 3093 754 | ERR309375 8 | |
| **DPP2 *.bam file run accession*** | ERR 2992 676 | ERR 2992 677 | ERR299267 5 | | ERR 2992 679 | ERR 2992 680 | ERR299267 8 | | ERR 2992 682 | ERR 2992 683 | ERR299268 1 | |
| **DPP3 *.bam file run accession*** | ERR 2985 803 | ERR 2985 804 | ERR298580 2 | | ERR 2985 806 | ERR 2985 807 | ERR298580 5 | | ERR 2985 809 | ERR 2985 810 | ERR298580 8 | |
| **Experiment Accession 100k*** | ERX 3156 15 | ERX 2986 848 | NA | NA | ERX 3156 16 | ERX 2986 847 | NA | NA | ERX 3156 17 | ERX 2986 849 | NA | NA |
| **Run Accession 100k *** | ERR 3090 788 | ERR 2984 276 | NA | NA | ERR 3090 789 | ERR 2984 275 | NA | NA | ERR 3090 790 | ERR 2984 277 | NA | NA |

153   * Using the Study Accession numbers in the European Nucleotide Archive all related data

154   files can be accessed, or accessed directly from

155   https://www.ebi.ac.uk/ena/data/view/accession, e.g.:

156   https://www.ebi.ac.uk/ena/data/view/PRJEB9846 (Study Accession Institution I),

157   https://www.ebi.ac.uk/ena/data/view/ERR972805 (Run Accession UKDD Institution I).

158

159   **Fig 1: Flowchart of the experimental design.** SP: sequence platform; DPP: data processing

160   pipeline

161

# Samples

163   All samples were obtained from outbreaks in commercial poultry holdings. Isolate

164   A/duck/England/36254/2014 was obtained from pooled intestinal material from index case

165   ducks (*Anas platyrhynchos domesticus*). Tissue homogenate material was inoculated into

166   embryonated chicken eggs and allantoic fluid was harvested at 1 day post-inoculation [14].

167   The Dutch isolate (A/chicken/Netherlands/EMC-3/2014) was obtained by passaging lung

168   material of a dead commercial layer hen *(Gallus gallus domesticus)* in MDCK cells twice and

169  harvesting the supernatant after approximately 40 hours post-inoculation [23]. The German

170  isolate (A/turkey/Germany/AR2485/2014) originated from lung tissue of a commercially kept

171  turkey (*Meleagris gallopavo)* and was passaged in embryonated chicken eggs [15]. (Table 1)


172

# Sequencing

## Institution I: SP1

175  RNA was extracted using a Qiagen QIAamp viral RNA mini kit (Qiagen, Germany)

176  according to the manufacturers' instructions except that carrier RNA was omitted from the

177  AVL lysis buffer and the sample was eluted in 50µl RNAse-free water. RNA was then

178  processed to double-stranded cDNA (cDNA Synthesis System, Roche) using random

179  hexamers and purified using magnetic beads (AmpureXP, Beckman Coulter, USA). The

180  double-stranded cDNA was diluted to 0.2 ng/µl and used to produce a sequencing library

181  using the NexteraXT kit (Illumina, USA). Libraries were then sequenced in paired-end mode

182  on an Illumina MiSeq (Illumina, USA), with run lengths varying from 2 x 75 bases (UKDD

183  virus) to 2 x 150 bases (NLCH and DETU viruses) depending on whether time-constraints

184  were implemented to provide a rapid response to an outbreak. Demultiplexing and removal of

185  sequencing adapters was done by the MiSeq RTA software to generate raw fastq files. SP1

186  process included a limited 12-cycle PCR enrichment of the library.  Post-hoc analysis showed

187  that duplication levels were less than 0.02% of the total reads which were considered to have

188  negligible impact on the results.


189

## Institution II: SP2

191  RNA was extracted using a combined approach with TRIzol (Thermo Fisher Scientific, USA)

192  and an RNeasy Kit (Qiagen, Germany). Further concentration and cleaning was done with

193  Agencourt RNAClean XP magnetic beads (Beckman Coulter, USA). RNA was quantified

194  using a Nanodrop UV spectrometer ND-1000 (Peqlab, Germany) and used as template for

195  cDNA synthesis with a cDNA Synthesis System (Roche, Germany) with random hexamers.

196  Fragmentation of the cDNA applying a target size of 300 bp was done with a Covaris M220

197  ultrasonicator. The sonicated cDNA was used for library preparation using Illumina indices

198  (Illumina, USA) on a SPRI-TE library system (Beckman Coulter, USA) using a SPRIworks

199  Fragment Library Cartridge II (for Roche FLX DNA sequencer; Beckman Coulter, USA)

200  without automatic size selection. Subsequently, upper and lower size exclusion of the library

201  was done with Ampure XP magnetic beads (Beckman Coulter, USA). The libraries were

202  quality checked using High Sensitivity DNA Chips and reagents on a Bioanalyzer 2100

203  (Agilent Technologies, Germany) and quantified via qPCR with a Kapa Library

204  Quantification Kit (Kapa Biosystems, USA) on a Bio-Rad CFX96 Real-Time System (Bio-

205  Rad Laboratories, USA). SP2 did not amplify sample nor library. Sequencing was done on an

206  Illumina MiSeq using MiSeq reagent kit v3 (Illumina, USA) resulting in paired end

207  sequences with a read length of 300. Demultiplexed and adapter-trimmed reads were used to

208  generate raw fastq files.

209

## Institution III: SP3

211  RNA was extracted using the High Pure RNA isolation kit (Roche Diagnostics, Germany)

212  according to manufacturer's instructions. RNA was converted to cDNA using the SuperScript

213  III Reverse Transcriptase kit (Invitrogen, Thermo Fisher, USA) as described previously [24],

214  and amplified by PCR using primers covering the full viral genome (S1 Table). All 32 PCR

215    fragments from approximately 400-600 nucleotides in length, were sequenced using the

216    454/Roche GS-FLX sequencing platform. The PCR fragments were pooled in equimolar ratio

217    and purified using the MinElute PCR Purification kit (Qiagen, Germany) according to the

218    manufacturer's instructions. Rapid Library preparation, Emulsion PCR and Next Generation

219    454-sequencing were performed according to instructions of the manufacturer (Roche

220    Diagnostics, Germany). Protocols are described in the following manuals: Rapid Library

221    Preparation Method Manual (Roche; May 2010), emPCR Amplification Method Manual –

222    Lib-L (Roche; May 2010) and Sequencing Method Manual (Roche; May 2010). All three

223    samples were sequenced in one run. Samples were pooled using MID adaptors to determine

224    which sequences came from which sample, each sample was assigned two different MID's.

225    Demultiplexing and basic trimming was done by CLC-bio software to generate raw fastq files

226    (S1 File).


227

## Data processing

### Institution I: DPP1

230    In the FluSeqID script (https://github.com/ellisrichardj/FluSeqID) the following steps are run

231    automatically: the mapping of raw sequence data to the host genome (BWA v0.7.12-r1039

232    [25]), extracting reads that do not map to the host (Samtools v1.2 [26]), assembling non-host

233    reads (Velvet v1.2.10 [27]), identification of the closest match for each genome segment

234    (BLAST v2.2.28 [28] using the custom databases generated from the Influenza Research

235    Database as indicated in the GitHub repository), mapping original data to the top reference

236    segments (BWA), calling new consensus sequences (vcf2consensus.pl), performing further

237    iterations of the last two steps to improve new consensus (IterMap), and finally outputting the

238    genome consensus sequence. The data processing pipeline has in-build defaults for k-mer and

239    coverage cut-off for de novo assembly, and the e-value cut-off for BLAST, which can be

240    changed via command line options (see https://github.com/ellisrichardj/FluSeqID). Since the

241    aligner (BWA-MEM) used  performs soft-clipping and ignores low quality data, quality

242    trimming is unnecessary. For mSNV analysis, the reads were mapped to the unified

243    consensus using BWA. Samtools was used to generate a pileup file which was then analysed

244    using custom python and R scripts to determine the depth of coverage and basecalls at each

245    position (available at https://github.com/ellisrichardj/MinorVar). The combination of BWA-

246    MEM and samtools was shown to be accurate for SNV identification [29]. In order to be

247    included in the final output the minimum basecall quality was 20 and the minimum mapping

248    quality was 50.

249

## Institution II: DPP2

251    Raw sequence data were analysed and mapped using the Genome Sequencer software suite

252    (v. 3.0; Roche, Mannheim, Germany) and the Geneious software suite (v. 9.0.5; Biomatters,

253    Auckland, New Zealand). Raw reads were trimmed and subsets of each trimmed dataset were

254    assembled *de novo* to generate reference sequences for each data set (Newbler Assembler of

255    Genome Sequencer software suite v. 3.0; Roche, Mannheim, Germany). The trimmed raw

256    influenza virus reads were mapped to the reference sequences (Newbler Mapper of Genome

257    Sequencer software suite v. 3.0; Roche, Mannheim, Germany). The output assemblies were

258    imported into the Geneious software suite (v. 9.0.5; Biomatters, Auckland, New Zealand) for

259    further analysis and processing. Regions of low and high coverage (threshold was 2 x

260    standard deviations from the mean for low and high coverage) and regions of low quality

261    (minimum quality/phred score 20) were evaluated and if necessary, excluded from further

262    analyses. Consensus sequences were generated and annotated using annotated reference

263    sequences. Sequences were compared, and annotations that matched with a similarity (>

264    90%) were copied. This was done on nucleotide sequences and also for translation in all six

265    reading frames. Annotations were manually inspected and curated. Trimmed raw reads of the

266    datasets or subsets thereof were mapped to the consensus, mapping was fine-tuned and

267    mSNVs were determined using generic SNP finder of the Geneious software suite, applying

268    parameters of maximum p-value of $10^{-5}$ and filter for strand bias. The threshold for SNP

269    identification was set at 1%, and variants were checked manually for accuracy.


## 270    Institution III: DPP3


271    Raw sequence data were analysed and mapped using the CLC Genomics software package,

272    workbench 8 (CLC Bio). Reads obtained by 454 sequencing were sorted by MID adaptor,

273    quality-trimmed, and analysed using the parameters as shown in S1 File. In short, after

274    sorting by MID, the sequence reads were trimmed at 30 nucleotides from the 3′ and 5′ ends to

275    remove all primer sequences. Data from the shared Illumina sequence files had already been

276    trimmed and were imported in CLC Bio without additional processing steps (S1 File). Reads

277    were initially aligned to their own reference sequences that were uploaded during the H5N8

278    outbreak (Gisaid accession numbers EPI-ISL-169282 (NLCH), EPI-ISL-167904 (UKDD)

279    and EPI-ISL-169273 (DETU)). Consensus sequences were automatically generated by CLC

280    after alignment to the reference, for detailed settings see S1 File. For the mSNV analysis the

281    raw data were mapped to the new custom-made consensus sequences per gene segment per

282    sample. Fastq files of these alignments were shared with the other institutions. The threshold

283    for mSNV identification was set at 1%, and registered minority variants were checked

284    manually for accuracy (minimal quality/phred score 20).


285

**Determining the influence of the DPP alignment steps versus DPPs mSNV**

**identification methods**

Data processing pipelines process raw data in several steps, roughly divided into trimming,

aligning data to a reference sequence, and variant calling (the mSNV identification

procedure). In order to determine to what extent the trimming and subsequent alignment

processes contributed to the observed differences the nucleotide coverage results obtained by

the three DPPs when aligning the same SP raw datasets were compared. To study the

influence of the mSNV identification process, quality-trimmed alignment files that had been

generated by each DPP and shared as *.bam files were subjected to the mSNV identification

process used in DPP3 to determine the differences in mSNV detection output when only the

alignment processes differed. DPP3 was randomly picked for this analyses, mSNV detection

parameters were set to the institutions default settings for mSNV identification using CLC-

bio software and can be seen in the S1 File.


# Data sharing

To test the applicability of real-time sequence data sharing within the COMPARE network,

all raw sequence data used in this study were uploaded to and shared via a "Data Hub" in the

environment of the European Nucleotide Archive (ENA). Each institution received its own

study accession in which all raw sequence data files and metadata files were assigned with

individual experimental accession numbers (Table 1). In addition to the sequence data, all

trimmed alignment files (*.bam) have been uploaded to the ENA. Using these hubs, sharing

between institutions was facilitated and immediate access to the data prior to the public

308  release was possible to enable joint evaluation and comparison. All data files have been made

309  publicly available via the ENA (https://www.ebi.ac.uk/ena).

310

## Designing the custom-made consensus sequences

312  Each institution produced a consensus sequence for the 8 influenza gene segments (PB2,

313  PB1, PA, HA, NP, NA, MP, NS) for each of the three viruses. The obtained consensus

314  sequences were aligned using the BioEdit sequence alignment editor (version 7.2.0) [30].

315  Raw sequence data from each SP were initially aligned to their own  reference sequences that

316  were uploaded during the H5N8 outbreak (Gisaid accession numbers EPI-ISL-

317  169282 (NLCH), EPI-ISL-167904 (UKDD) and xxx (DETU)).

318

## mSNV analysis comparison

320  For the mSNV analyses the custom-made consensus for each virus isolate was used as a

321  reference for mapping, thereby standardizing positions within the genome to make

322  comparison between institutions easier. To avoid unnecessary increases in analytical time and

323  memory, datasets were down-sampled to 100.000 reads per sample when needed. Each DPP

324  produced a report on the identified mSNVs in a tabulated format. The analysis output files

325  were filtered for mSNVs only, thereby ignoring detected nucleotide insertions and deletions

326  (InDels). There is a current lack of data or evidence-based approaches on how to calculate the

327  required sequence depth (i.e. coverage) for mSNV analyses an evidence-based. In this study,

328  a minimum coverage threshold for the identification of mSNVs was applied. This minimum

329  nucleotide coverage (i.e. number of reads per nucleotide after trimming) was determined

330 using a basic sample size calculation method, $n = \log \beta / \log p'$ [31]. Here $\beta$ represents the

331 required power (e.g. for 95% chance of detecting something $\beta = 0.05$), and $p'$ is 1 - the

332 proportion of events that you want to detect. For a 95% certainty of detecting a variant at 1%,

333 a minimum coverage of 298 reads per position is needed. For variants that occur in $\geq$5% of

334 reads, the number of reads required is >58, and for variants that occur in $\geq$10% of the reads

335 the minimum coverage is >28. For the mSNV identification literature commonly uses the

336 mSNV cut-off frequencies of $\geq$10%, $\geq$5% and $\geq$1%. However, it needs to be noted that these

337 cut-off values are arbitrary. Therefore, where depth of coverage was sufficient, this study will

338 report mSNV detected with a frequency of $\geq$1%, but initial comparisons started with

339 positions showing mSNVs with frequencies of $\geq$10% in at least one of the SP/DPP

340 combinations, followed by those with mSNV of $\geq$5% -<10%, and lastly those $\geq$1%-<5%. For

341 all those positions identified, the number of reads and number of variant nucleotides in all

342 other SP/DPP combinations for that position will be noted regardless of frequencies.

343

# Results

345 In order to determine the comparability of consensus sequences and mSNV identification the

346 biological samples, the sequence data from 3 SPs and the *.bam quality-trimmed alignment

347 files of raw data of 3 influenza A/H5N8 viruses were shared. All data sets were subsequently

348 analysed in 3 different DPPs. The resulting 9 mSNV reports per virus (3 SP data sets each

349 analysed in 3 DPPs) were evaluated for comparability of mSNV identification using different

350 combinations of SP and DPP.

351

## Data sharing

Data sharing using the COMPARE "Data Hub" provided by ENA proved to be easy, quick and successful. The "Data Hub" enables the File Transport Protocol (FTP) protected upload and download of large data files and facilitates sharing between collaborators with the possibility to evaluate and compare all data prior to their public release by generating and specifically sharing accession numbers using standard ENA procedures. The Data Hub used an influenza virus sample checklist. In addition, data sets are ultimately made publicly and through the INSDC network globally available and accessible in real-time as required without further upload to a different repository. Full details of the COMPARE Data Hub system are available in a submitted manuscript [22]. In summary, this process was suitable for quick data sharing in an outbreak scenario.

## Designing the custom-made consensus sequences

For each of the 8 gene segments of the 3 viruses separately, 9 initial consensus sequences (3 SPs x 3 DPPs) were generated, resulting in 72 consensus sequences per virus. The custom-made consensus sequence per virus and per gene segment was 1) trimmed to a length represented by all 9 initial consensus sequences and 2) nucleotides had to be identical to at least 6/9 consensus sequences to be included. Although some sequences contained insertions or deletions, those could always be corrected for using the other SP sequences following the criteria mentioned previously. This resulted in a unique custom-made consensus for each gene segment for all three viruses.

# Consensus sequences

When ignoring insertions and deletions in the homopolymer regions of the 454 data for most
gene segments the identified consensus sequences were identical regardless of the SP and
DPP combinations used with the exemption of the differences mentioned in Table 2.
However, the number of insertions and deletions in homopolymer regions of the SP3
sequences were considerable in all 3 viruses. There was no clear difference in the number of
insertions and deletions related to homopolymer regions between the different DPPs (20, 17
and 18 for DPP 1, 2 and 3 respectively). Nucleotide differences that were not related to
homopolymer regions were only observed for sequences obtained in SP3 and SP2 data when
processed in DPP1.

**Table 2. The differences in consensus sequences obtained from each SP/DPP**

**combination, sorted per virus and per gene segment.**

| Virus | Segment | Start* | End | Number of InDels at homopolymer regions** | Other nucleotide differences*** |
|-------|---------|--------|-----|-------------------------------------------|--------------------------------|
| NLCH | PB2 | 1 | 2280 | 2 (DPP1)<br>2 (DPP3) | C506A (SP3)<br>G2101R (SP3) |
| | PB1 | 1 | 2277 | 1 (DPP1/DPP2/DPP3)<br>1 (DPP1/DPP2)<br>1 (DPP2/DPP3)<br>1 (DPP3) | A949W (SP3)<br>2272 ins AAG (SP2) |
| | PA | -6# | 2190 | 1 (DPP1/DPP2)<br>2 (DPP1) | ND |
| | HA | 7 | 1704 | 1 (DPP2/DPP3) | A427W (SP2) |
| | NP | 1 | 1497 | 1 (DPP1) | C420Y (SP3) |
| | NA | 4 | 1419 | ND | ND |
| | MP | -5# | 982 | ND | ND |
| | NS | 1 | 838 | ND | ND |
| DETU | PB2 | 1 | 2287 | 1 (DPP1/DPP2/DPP3)<br>3 (DPP1) | 2272 Del A (SP3) |
| | PB1 | 1 | 2277 | 1 (DPP1/DPP2/DPP3)<br>1 (DPP1) | T956C (SP3) |

| | | | | | |
|---|---|---|---|---|---|
| | | | | 1 (DPP2) | |
| | | | | 1 (DPP3) | |
| | PA | 7 | 2189 | 1 (DPP1/DPP2) | ND |
| | HA | 1 | 1728 | 1 (DPP2/DPP3) | ND |
| | NP | 1 | 1497 | 2 (DPP3) | ND |
| | NA | 1 | 1413 | 1 (DPP1) | 778 ins CCA (SP3) |
| | MP | -1# | 982 | 1 (DPP2) | ND |
| | NS | 2 | 838 | ND | ND |
| **UKDD** | PB2 | 1 | 2298 | 1 (DPP1/DPP2/DPP3) | C504T (SP3) |
| | | | | 1 (DPP3) | C506M (SP3) |
| | PB1 | 1 | 2277 | 1 (DPP1/DPP2/DPP3) | T951W (SP3) |
| | | | | 1 (DPP2/DPP3) | |
| | PA | 1 | 2151 | 2 (DPP1) | ND |
| | | | | 1 (DPP2) | |
| | HA | 1 | 1704 | 1 (DPP2/DPP3) | ND |
| | NP | 1 | 1497 | 1 (DPP3) | T1003Y (SP2) |
| | NA | 4 | 1420 | ND | 782 del TA (SP3) |
| | MP | -5# | 982 | 1 (DPP2) | ND |
| | NS | -5# | 849 | ND | ND |

387    The letter in brackets represents the DPP (column 5) or the SP (column 6) where the

388    insertions/deletions or mutations were detected. InDel: insertions or deletion; SP: Sequence

389    platform; DPP: Data processing pipeline; ND: not detected. * Start is counted from the ATG

390    start codon; ** Exclusively identified in SP3 sequence data, InDels related to homopolymer

391    regions; *** Exclusively identified in DPP1; # '-' indicates the number of nucleotides before

392    the ATG start codon included in the consensus

393

394

395    In summary, the homopolymer errors inherent in the 454 dataset caused problems for all

396    DPPs, as expected. Consensus sequences generated by DPP1 from SP3 (454) data showed

397    some unexpected differences, but it performed well with the SP1 data formats it was designed

398    for and reasonably well with SP2 data. Overall, the consensus sequences can be reproduced

399    by all DPPs using Illumina data but that the analysis of the 454 data from SP3 was more

400  problematic, as it would require editing of the sequences at homopolymer regions. Consensus

401  sequences from this study can be found in the S2 Table.

402

## The mSNV analysis comparison

### Nucleotide coverage and the influence of DPP-dependent alignment

405  The observed number of reads per nucleotide (referred to as nucleotide coverage) differed

406  depending on the SP/DPP combination. All DPPs handled both 454 and Illumina data

407  formats, although some modifications (settings for the bwa mapper to handle single end 454

408  data) were required for DPP1, which was specifically designed for Illumina paired-end reads.

409  The observed nucleotide coverages showed near to identical profiles for all three viruses. The

410  coverage results obtained from the three different SPs and DPPs for the NLCH virus (Fig 2)

411  and for the other two viruses (S1 Figure) were plotted. In general, lower nucleotide coverage

412  was observed at the termini of each gene segment. The SP3 data showed more variation in

413  nucleotide coverage within gene segments compared to SP1 and SP2 data, due to the

414  sequencing of 32 PCR amplicons. The non-normalised number of raw sequence reads and

415  influenza virus reads per virus per SP can be found in the S3 Table.

416

417  **Fig 2: Nucleotide coverage.** The non-normalised nucleotide coverage displayed as number

418  of nucleotides per position for full genome sequences of the NLCH virus reads mapped to the

419  NLCH reference sequences. Panel A shows the coverage results for the same SP dataset in

420  the three different DPPs (DPP1: purple; DPP2: orange; DPP3 grey) for each of the SP

421  datasets. Panel B shows the coverage when the same DPP is used to analyse data from the

422 three different SPs (SP1: lilac; SP2: yellow; SP3: green) for each of the DPPs. The X-axis

423 represents the position in the genome, the Y-axis represents the number of sequence reads per

424 position.

425

426 The differences in nucleotide coverage were visualized for the three different SP raw datasets

427 analysed with the same DPP (Fig 2A). Overall, SP3 data (green lines) showed a lower

428 coverage compared to SP1 (purple) and SP2 data (yellow). The overall coverage for SP1 and

429 SP2 data was similar with small variations for different viruses and DPPs. The shorter read

430 lengths in SP1 virus data did not appear to have influenced the overall nucleotide coverage

431 substantially.

432 The differences in nucleotide coverage introduced by different alignment procedures were

433 also assessed, by comparing the coverage results for each SP raw dataset analysed with the

434 three different DPPs (Fig 2B). DPP2 (orange lines) generally retained the highest nucleotide

435 coverage for data from the different SPs. However, DPP3 (grey lines) generally also retained

436 high coverage for SP3 data, for which it was optimized. The nucleotide coverage of SP3 data

437 showed larger variation between the three different DPPs, leading to differences in nucleotide

438 coverage up to 50% depending on the DPP, because DPP1 and DPP2 were not optimized for

439 this SP. Data from SP2 were handled very similar by all three DPPs.

440 In conclusion, both the SP and the DPP influenced the number of reads per nucleotide

441 position. SP3 showed the lowest output in number of reads compared to SP1 and SP2

442 Illumina data. The influence of the DPP depended highly on the data input, with best DPP

443 performance for the SP dataset for which it was optimized.

444

445 **The mSNV identification**

446   The mSNV identification thresholds were set to ≥1% in all DPPs. Because of the high

447   number of mSNVs identified, the comparison of these mSNVs started with a manually set

448   arbitrary threshold of ≥10% that was subsequently decreased to ≥5%, and ≥1%. A mSNV

449   position was identified when at least 1 of the SP/DPP combinations showed a variant that

450   exceeded the frequency threshold, and when the coverage at that position exceeded the

451   minimum number of reads needed to detect that variant with a 95% probability, as described

452   previously. The presence of mSNV and coverage for all SP/DPP combinations were

453   compared for each of the positions in which a mSNV had been detected in at least one of the

454   combinations. The coverages indicated for those positions where no mSNVs were detected

455   were derived from the alignment files and were not subjected to possible additional read

456   filtering parameters in the mSNV identification process. The average quality (Q-score/phred

457   score) was set to or exceeding 20.

458   Ten positions across the three virus genomes were identified with mSNVs occurring in ≥10%

459   of reads. Three of the mSNVs (NLCH:PB2 G1879A , NLCH:PB2 G2101A and DETU:HA

460   T963C) were detected in all SP/DPP combinations but with slightly different relative

461   abundance. The other mSNVs were identified in only one (n=6) or two (n=1) of the SP/DPP

462   combinations (Table 3).

463

464   **Table 3. The minority variants occurring in at least one of the sequence platform - data**

465   **processing pipelines as a ≥5% variant.**

| Virus | Position | Sequence platform | Data processing pipeline | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1 | | 2 | | 3 | |
| | | | Minor variants | Percentage | Minor variants | Percentage | Minor variants | Percentage |
| NLCH | PB2.1879 G→A | 1 | 81/1301 | 6,2% | 246/2716 | 9,1% | 112/1203 | 9,3% |
| | | 2 | 47/956 | 4,9% | 117/1137 | 10,3% | 114/1064 | 10,7% |
| | | 3 | 49/530 | 9,2% | 131/1341 | 9,8% | 129/1338 | 9,6% |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | **PB2.2101 G→A** | 1 | 53/1118 | 4,7% | 261/2704 | 9,7% | 110/897 | 12,3% |
| | | 2 | 21/1578 | 1,3% | 125/1875 | 6,7% | 121/1463 | 8,3% |
| | | 3 | 13/542 | 2,4% | 199/1433 | 13,9% | 199/1435 | 13,9% |
| | **PB2.2277 T→G** | 1 | ND/479 | <1% | 86/1008 | 8,5% | 33/190 | 17,4% |
| | | 2 | ND/557 | <1% | ND/623 | <1% | ND/534 | <1% |
| | | 3 | ND/680 | <1% | ND/1117 | <1% | ND/1024 | <1% |
| | **PB1.87 A→G** | 1 | ND/818 | <1% | ND/1754 | <1% | ND/1114 | <1% |
| | | 2 | 25/230 | 10,9% | ND/376 | <1% | ND/328 | <1% |
| | | 3 | ND/275 | <1% | ND/537 | <1% | ND/537 | <1% |
| | **PB1.2240 G→C** | 1 | ND/664 | <1% | 54/1341 | 4,0% | 38/418 | 9,1% |
| | | 2 | ND/1231 | <1% | ND/1271 | <1% | ND/1233 | <1% |
| | | 3 | ND/161 | <1% | ND/277 | <1% | ND/276 | <1% |
| | **PB1.2268 A→G** | 1 | ND/336 | <1% | 29/641 | 4,5% | 11/176 | 6,3% |
| | | 2 | ND/993 | <1% | ND/1026 | <1% | ND/1002 | <1% |
| | | 3 | ND/53 | <1% | ND/159 | <1% | ND/148 | <1% |
| | **HA.104 A→G** | 1 | ND/733 | <1% | ND/1761 | <1% | ND/1151 | <1% |
| | | 2 | ND/437 | <1% | ND/1370 | <1% | ND/1156 | <1% |
| | | 3 | ND/1 | <1% | ND/105 | <1% | 12/105 | 11,4% |
| | **HA.1689 T→C** | 1 | ND/390 | <1% | ND/694 | <1% | 11/217 | 5,1% |
| | | 2 | ND/2018 | <1% | ND/4083 | <1% | ND/3979 | <1% |
| | | 3 | ND/937 | <1% | ND/1669 | <1% | ND/1680 | <1% |
| | **NP.105 A→G** | 1 | ND/182 | <1% | ND/449 | <1% | ND/343 | <1% |
| | | 2 | 83/1507 | 5,5% | ND/1890 | <1% | ND/1804 | <1% |
| | | 3 | ND/89 | <1% | ND/704 | <1% | ND/702 | <1% |
| | **NP.1239 A→T** | 1 | 32/2428 | 1,3% | 279/5410 | 5,2% | ND/3092 | <1% |
| | | 2 | ND/2345 | <1% | ND/2643 | <1% | ND/2453 | <1% |
| | | 3 | ND/1711 | <1% | ND/2111 | <1% | ND/2117 | <1% |
| | **NP.1489 G→A** | 1 | ND/182 | <1% | 26/336 | 7,7% | ND/172 | <1% |
| | | 2 | ND/436 | <1% | ND/452 | <1% | ND/444 | <1% |
| | | 3 | ND/1320 | <1% | ND/1799 | <1% | ND/1799 | <1% |
| | **NS.833 A→T** | 1 | ND/187 | <1% | ND/287 | <1% | 5/88 | 5,7% |
| | | 2 | ND/1224 | <1% | ND/1327 | <1% | ND/1284 | <1% |
| | | 3 | ND/1367 | <1% | ND/2430 | <1% | ND/2333 | <1% |
| **DETU** | **PB2.1054 T→C** | 1 | 69/1369 | 5,0% | 168/2637 | 6,4% | 97/1304 | 7,4% |
| | | 2 | 60/1477 | 4,1% | 115/1836 | 6,3% | 99/1605 | 6,2% |
| | | 3 | 6/392 | 1,5% | 94/2038 | 4,6% | 48/1054 | 4,6% |
| | **PB2.2257 A→C** | 1 | ND/867 | <1% | ND/1563 | <1% | 24/463 | 5,2% |
| | | 2 | ND/531 | <1% | ND/581 | <1% | ND/378 | <1% |
| | | 3 | ND/893 | <1% | ND/2286 | <1% | ND/1346 | <1% |
| | **PB2.2277 T→G** | 1 | ND/644 | <1% | 52/1150 | 4,5% | 27/307 | 8,8% |
| | | 2 | ND/418 | <1% | ND/472 | <1% | ND/284 | <1% |
| | | 3 | ND/1208 | <1% | ND/1948 | <1% | ND/1209 | <1% |

| Mutation | | | | | | | |
|---|---|---|---|---|---|---|---|
| PB1.14 C→T | 1 | ND/144 | <1% | 48/433 | 11,1% | ND/239 | <1% |
| | 2 | ND/90 | <1% | ND/355 | <1% | ND/304 | <1% |
| | 3 | ND/562 | <1% | ND/792 | <1% | ND/496 | <1% |
| PB1.23 T→G | 1 | ND/207 | <1% | 30/535 | 5,6% | ND/315 | <1% |
| | 2 | ND/103 | <1% | ND/365 | <1% | ND/319 | <1% |
| | 3 | ND/699 | <1% | ND/950 | <1% | ND/609 | <1% |
| PB1.87 A→G | 1 | ND/744 | <1% | ND/1644 | <1% | ND/1076 | <1% |
| | 2 | 49/365 | 13,4% | ND/677 | <1% | ND/576 | <1% |
| | 3 | ND/721 | <1% | ND/1156 | <1% | ND/793 | <1% |
| PB1.2240 G→C | 1 | ND/757 | <1% | 23/1517 | 1,5% | 26/515 | 5,0% |
| | 2 | ND/944 | <1% | ND/985 | <1% | ND/806 | <1% |
| | 3 | ND/274 | <1% | ND/439 | <1% | ND/253 | <1% |
| PB1.2268 A→G | 1 | 5/470 | 1,1% | 33/928 | 3,6% | 22/278 | 7,9% |
| | 2 | ND/798 | <1% | ND/829 | <1% | ND/671 | <1% |
| | 3 | ND/109 | <1% | ND/259 | <1% | ND/123 | <1% |
| PB1.2271 A→G | 1 | 12/446 | 2,7% | 59/901 | 6,5% | 16/263 | 6,1% |
| | 2 | ND/729 | <1% | 47/810 | 5,8% | 40/649 | 6,2% |
| | 3 | 1/32 | 3,1% | ND/123 | <1% | 2/83 | 2,4% |
| HA.867 C→T | 1 | 59/1533 | 3,8% | 206/3183 | 6,5% | 104/1537 | 6,8% |
| | 2 | 59/2031 | 2,9% | 150/2525 | 5,9% | 127/2253 | 5,6% |
| | 3 | 11/180 | 6,1% | 48/647 | 7,4% | 28/385 | 7,3% |
| HA.963 T→C | 1 | 122/1401 | 8,7% | 446/3071 | 14,5% | 189/1419 | 13,3% |
| | 2 | 90/1517 | 5,9% | 318/2189 | 14,5% | 247/1828 | 13,5% |
| | 3 | 5/69 | 7,2% | 107/606 | 17,7% | 47/293 | 16,0% |
| NP.1491 C→A | 1 | ND/278 | <1% | 71/583 | 12,2% | ND/206 | <1% |
| | 2 | ND/723 | <1% | ND/769 | <1% | ND/692 | <1% |
| | 3 | ND/799 | <1% | ND/2031 | <1% | ND/1206 | <1% |
| NA.65 T→C | 1 | 19/503 | 3,8% | 52/1229 | 4,2% | 16/467 | 3,4% |
| | 2 | 20/662 | 3,0% | 50/1104 | 4,5% | 45/992 | 4,5% |
| | 3 | 24/557 | 4,3% | 53/1099 | 4,8% | 37/727 | 5,1% |
| NA.78 T→C | 1 | 23/599 | 3,8% | 57/1403 | 4,1% | 20/557 | 3,6% |
| | 2 | 21/692 | 3,0% | 55/1147 | 4,8% | 50/1033 | 4,8% |
| | 3 | 23/580 | 4,0% | 51/1124 | 4,5% | 37/735 | 5,0% |
| NA.89 T→C | 1 | 23/713 | 3,2% | 55/1670 | 3,3% | 22/651 | 3,4% |
| | 2 | 23/798 | 2,9% | 56/1261 | 4,4% | 50/1134 | 4,4% |
| | 3 | 24/580 | 4,1% | 55/1196 | 4,6% | 40/775 | 5,2% |
| NA.117 T→C | 1 | 37/908 | 4,1% | 87/2140 | 4,1% | 36/818 | 4,4% |
| | 2 | 28/1102 | 2,5% | 67/1631 | 4,1% | ND/1459 | <1% |
| | 3 | 22/531 | 4,1% | 57/1276 | 4,5% | 42/812 | 5,2% |
| NA.126 T→C | 1 | 37/983 | 3,8% | 83/2294 | 3,6% | 36/876 | 4,1% |
| | 2 | 31/1126 | 2,8% | 72/1676 | 4,3% | 65/1502 | 4,3% |
| | 3 | 26/519 | 5,0% | 62/1395 | 4,4% | 43/812 | 5,3% |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| UKDD | PB2.2277 T→G | 1 | ND/415 | <1% | 28/507 | 5,5% | ND/475 | <1% |
| | | 2 | ND/589 | <1% | ND/620 | <1% | ND/601 | <1% |
| | | 3 | ND/1140 | <1% | ND/1996 | <1% | ND/2065 | <1% |
| | PB1.87 A→G | 1 | ND/387 | <1% | ND/440 | <1% | ND/439 | <1% |
| | | 2 | 26/327 | 8,0% | 32/395 | 8,1% | ND/351 | <1% |
| | | 3 | ND/617 | <1% | ND/1133 | <1% | ND/1136 | <1% |
| | PB1.728 C→A | 1 | ND/750 | <1% | ND/832 | <1% | ND/836 | <1% |
| | | 2 | ND/776 | <1% | 52/928 | 5,6% | ND/829 | <1% |
| | | 3 | ND/2459 | <1% | ND/4290 | <1% | ND/4293 | <1% |
| | PB1.730 C→T | 1 | ND/742 | <1% | ND/824 | <1% | ND/826 | <1% |
| | | 2 | ND/767 | <1% | 57/1008 | 5,7% | ND/832 | <1% |
| | | 3 | ND/2339 | <1% | ND//4286 | <1% | ND/4289 | <1% |
| | PB1.883 G→C | 1 | ND/942 | <1% | ND/997 | <1% | ND/997 | <1% |
| | | 2 | ND/1689 | <1% | ND/1865 | <1% | ND/1760 | <1% |
| | | 3 | ND/2479 | <1% | 47/690 | 6,8% | ND/3681 | <1% |
| | PA.49 G→C | 1 | ND/103 | <1% | 6/117 | 5,1% | ND/115 | <1% |
| | | 2 | ND/337 | <1% | ND/435 | <1% | ND/392 | <1% |
| | | 3 | ND/111 | <1% | ND/207 | <1% | ND/204 | <1% |
| | PA.82 C→T | 1 | ND/155 | <1% | ND/180 | <1% | ND/177 | <1% |
| | | 2 | ND/695 | <1% | ND/809 | <1% | ND/745 | <1% |
| | | 3 | ND/64 | <1% | ND/247 | <1% | 30/248 | 12,1% |
| | NS.811 G→T | 1 | ND/221 | <1% | 17/270 | 6,3% | ND/249 | <1% |
| | | 2 | ND/2452 | <1% | ND/2725 | <1% | ND/2557 | <1% |
| | | 3 | ND/3117 | <1% | ND/4125 | <1% | ND/4139 | <1% |

Colours display the variant frequency with ≥10% (green), 5-10% (purple) and <5% (pink).

ND: not detected.




Thirty-seven positions were identified with mSNVs occurring in ≥5% of reads. Of those, the same mSNV was identified in all SP/DPP combinations for 9 positions (24,3%), in seven or eight of the SP/DPP combinations for 2 positions (5,4%) and in at least two SP/DPP combinations for 19 positions (51.4%), although not always in a frequency of ≥5%. However, for 18 positions (48.6%) the mSNV was not reproduced at a ≥1% frequency in any of the other SP/DPP combinations (Table 3). Focussing on the separate SP data analysed in the 3

476    DPPs, most of the identified positions with ≥5% mSNVs in at least 1 SP/DPP combination

477    were identified in SP1 data (47%) followed by SP2 (29%) and SP3 (24%) data.

478    Looking at the ≥5% mSNV reproducibility per SP dataset in all three DPPs within these

479    thirty-seven positions, forty-eight SP datasets showed a ≥5% mSNV in at least one of the

480    DPP outputs. Additionally, for eleven positions, all in the DETU virus, the variant was

481    reproduced by all DPPs, however at a <5% frequency (for instance SP3 data at PB2.1054,

482    and SP1 and SP2 data at NA.65) In 53% (31/59) of cases the same mSNVs from 1 SP dataset

483    was reproduced in all three DPP's in at least a ≥1% frequency, in 31% (18/59) of cases the

484    variant was only detected in 1 DPP even though coverage in the other DPPs was theoretically

485    high enough to detect variants at a 1% level.

486    Lowering the threshold value to a mSNV frequency of ≥1% resulted in a large increase in the

487    number of positions identified with mSNVs. To investigate the reproducibility of these

488    mSNVs, the data for all 3 viruses was combined per SP in the three DPPs (influence of DPP),

489    and per DPP analysing data from the three SPs (influence of SP). The genome positions with

490    ≥1% variants were listed per SP/DPP combination and entered in the program Venny 2.1 that

491    calculated the overlapping positions as a fraction of the total number of positions between the

492    SP/DPP combinations as compared to the total number of positions, resulting in Fig 3. It

493    needs to be noted that especially SP3 did not always reach the minimum coverage

494    requirements and may therefore not be suitable to detect low-frequency variants with (see

495    also table 4). Positions where the coverage in one or more of the nine SP/DPP combinations

496    didn't meet the minimum required coverage of 298 were not included in the comparison in

497    Fig 3. The reproducibility of ≥1% variants using one SP dataset in all three DPPs was 10%,

498    9.4% and 31.1% for SP1, SP2 and SP3 sequences, respectively. The reproducibility of ≥1%

499    variants using raw data of a virus sequenced in three different SPs was 20%, 7.4% and 22.6%

500    for DPP1, DPP2 and DPP3 respectively (Fig 3). Most ≥1% variants were not reproduced by

501    any of the other DPPs processing the same SP data (~75%) for SP1 and SP2 data. This was

502    less for SP3 data but this might be due to the fact that many positions identified in SP3 data

503    did not meet the minimum coverage criteria and were therefore discarded.

504    **Fig 3**: The reproducibility of ≥1% variants with sufficient coverage (>298) for all sequence

505    data combined. Each figure shows the number of ≥1% variants detected per sequence

506    platform (SP, top row) and data processing pipeline (DPP, bottom row) for SP1/DPP1 (left

507    column), SP2/DPP2 (middle column), and SP3/DPP3 (right column). The colours represent

508    the different DPPs and SPs respectively, in which the >1% variants were detected: SP1/DPP1

509    (purple), SP2/DPP2 (yellow) and SP3/DPP3 (green). Positions with ≥1% variants that were

510    identified in more than one of the SPs or DPPs respectively are displayed in the overlapping

511    coloured areas, the centre part representing the number of ≥1% variants that were detected

512    with all three DPPs (top row) or SPs (bottom row). The total number of positions with ≥1%

513    variants detected was 250in SP1, 213 in SP2, 45 in SP3, and 50 in DPP1, 353 in SP2, and 93

514    in SP3. This figure was produced using Venny 2.1.

515

516    For brevity, the detailed results for the HA gene segment of the DETU virus are shown in

517    Table 4. This virus segment was chosen because it showed the best reproducibility of results

518    for ≥5% minority variants in all SP/DPP combinations. In the DETU HA segment, 33

519    positions containing a mSNV occurring in ≥1% of reads with sufficient coverage (≥298

520    reads) were identified. Only 3 of these positions (9%) were identified in all SP/DPP

521    combinations. The majority of the positions (25/33, 76%) were only identified in one of the

522    nine SP/DPP combinations. However, it needs to be noted that the SP3 data coverage was

523    insufficient in all three DPPs to detect ≥1% variants for 11 of those positions (Table 4).

524    Although a comparison between the frequencies of the detected mSNVs might be

525    appropriate, based on these results where even absence vs. presence of the mSNVs is poorly

526    comparable further in-depth analyses on these frequencies is not performed because of its

527    limited value.

528    **Table 4. The minority variants occurring in at least one of the sequence platform - data**

529    **processing pipelines as a ≥1% variant in the HA segment of the DETU sample with a**

530    **minimum coverage of 298 reads at that position.**

| Position | Sequence platform | Data processing pipeline | | | | | |
|---|---|---|---|---|---|---|---|
| | | **1** | | **2** | | **3** | |
| | | **Minor variants** | **Percentage** | Minor variants | **Percentage** | **Minor variants** | **Percentage** |
| **HA.170 T→A** | 1 | ND/935 | <1% | ND/2191 | <1% | ND/1348 | <1% |
| | 2 | ND/300 | <1% | 11/693 | 1,59% | ND/551 | <1% |
| | 3 | ND/82* | <1%* | ND/245* | <1%* | ND/210* | <1%* |
| **HA.170 T→C** | 1 | ND/935 | <1% | ND/2191 | <1% | ND/1348 | <1% |
| | 2 | ND/300 | <1% | 18/693 | 2,60% | ND/551 | <1% |
| | 3 | ND/82* | <1%* | ND/245* | <1%* | ND/210* | <1%* |
| **HA.171 C→A** | 1 | ND/931 | <1% | ND/2184 | <1% | ND/1339 | <1% |
| | 2 | ND/323 | <1% | 12/698 | 1,72% | ND/558 | <1% |
| | 3 | ND/82* | <1%* | ND/245* | <1%* | ND/210* | <1%* |
| **HA.194 C→A** | 1 | ND/991 | <1% | ND/2397 | <1% | ND/1455 | <1% |
| | 2 | ND/353 | <1% | 22/701 | 3,14% | ND/553 | <1% |
| | 3 | ND/58* | <1%* | ND/250* | <1%* | ND/212* | <1%* |
| **HA.195 C→A** | 1 | ND/995 | <1% | ND/2390 | <1% | ND/1464 | <1% |
| | 2 | ND/356 | <1% | 20/701 | 2,85% | ND/553 | <1% |
| | 3 | ND/55* | <1%* | ND/250* | <1%* | ND/212* | <1%* |
| **HA.268 C→T** | 1 | ND/1140 | <1% | ND/2580 | <1% | ND/1626 | <1% |
| | 2 | ND/1293 | <1% | 25/1563 | 1,60% | ND/1338 | <1% |
| | 3 | ND/88* | <1%* | ND/252* | <1%* | ND/212* | <1%* |
| **HA.272 A→T** | 1 | ND/1156 | <1% | ND/2593 | <1% | ND/1639 | <1% |
| | 2 | 17/1424 | 1,19% | 20/1563 | 1,28% | ND/1404 | <1% |
| | 3 | ND/81* | <1%* | ND/253* | <1%* | ND/213* | <1%* |
| **HA.407 G→T** | 1 | ND/1144 | <1% | ND/2364 | <1% | ND/1553 | <1% |
| | 2 | ND/1773 | <1% | 31/2121 | 1,46% | ND/1855 | <1% |
| | 3 | ND/74* | <1%* | ND/237* | <1%* | ND/212* | <1%* |
| **HA.407 G→A** | 1 | ND/1144 | <1% | 27/2364 | 1,14% | ND/1553 | <1% |
| | 2 | ND/1773 | <1% | ND/2121 | <1% | ND/1856 | <1% |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | 3 | ND/74* | <1%* | ND/237* | <1%* | ND/212* | <1%* |
| **HA.418 A→G** | 1 | ND/1111 | <1% | ND/2319 | <1% | ND/1492 | <1% |
| | 2 | 29/2195 | 1,32% | 38/2513 | 1,51% | ND/2197 | <1% |
| | 3 | ND/69* | <1%* | ND/237* | <1%* | ND/212* | <1%* |
| **HA.453 T→G** | 1 | ND/1339 | <1% | 29/2736 | 1,06% | ND/1811 | <1% |
| | 2 | ND/2342 | <1% | ND/2695 | <1% | ND/2384 | <1% |
| | 3 | ND/91* | <1%* | ND/193* | <1%* | ND/179* | <1%* |
| **HA.560 A→G** | 1 | 43/1587 | 2,71% | 113/3385 | 3,34% | 55/1517 | 3,63% |
| | 2 | 56/2397 | 2,34% | 145/2912 | 4,98% | 113/2495 | 4,53% |
| | 3 | 21/884 | 2,38% | 72/1754 | 4,10% | 43/1245 | 3,45% |
| **HA.715 C→T** | 1 | ND/1663 | <1% | 62/3832 | 1,62% | 24/1582 | 1,52% |
| | 2 | 26/2283 | 1,14% | 55/2722 | 2,02% | 50/2420 | 2,07% |
| | 3 | ND/531 | <1% | 20/1883 | 1,06% | 15/1245 | 1,20% |
| **HA.867 C→T** | 1 | 59/1533 | 3,85% | 206/3183 | 6,47% | 104/1537 | 6,77% |
| | 2 | 59/2031 | 2,90% | 150/2525 | 5,94% | 127/2253 | 5,64% |
| | 3 | 11/180 | 6,11% | 48/647 | 7,42% | 28/385 | 7,27% |
| **HA.963 T→C** | 1 | 122/1401 | 8,71% | 446/3071 | 14,52% | 189/1419 | 13,32% |
| | 2 | 90/1517 | 5,93% | 318/2189 | 14,53% | 247/1828 | 13,51% |
| | 3 | 5/69 | 7,25% | 107/606 | 17,66% | 47/293 | 16,04% |
| **HA.100 0 A→C** | 1 | ND/1409 | <1% | 48/2962 | 1,62% | ND/1873 | <1% |
| | 2 | ND/1629 | <1% | ND/1919 | <1% | ND/1645 | <1% |
| | 3 | ND/84* | <1%* | ND/614 | <1% | ND/293* | <1%* |
| **HA.117 7 G→A** | 1 | ND/1222 | <1% | ND/2224 | <1% | ND/1597 | <1% |
| | 2 | ND/1652 | <1% | 34/1901 | 1,79% | ND/1724 | <1% |
| | 3 | ND/289* | <1%* | ND/549 | <1% | ND/270 | <1% |
| **HA.118 3 A→G** | 1 | ND/1210 | <1% | ND/2226 | <1% | ND/1589 | <1% |
| | 2 | ND/1770 | <1% | ND/1892 | <1% | ND/1723 | <1% |
| | 3 | ND/280* | <1%* | 6/547 | 1,10% | ND/268* | <1%* |
| **HA.119 9 T→G** | 1 | ND/1182 | <1% | ND/2124 | <1% | ND/1518 | <1% |
| | 2 | ND/1615 | <1% | 27/1899 | 1,42% | ND/1732 | <1% |
| | 3 | ND/296* | <1%* | ND/545 | | ND/266* | <1%* |
| **HA.126 3 A→G** | 1 | 16/963 | 1,66% | 57/1841 | 3,10% | 26/954 | 2,73% |
| | 2 | 26/1924 | 1,35% | 56/2207 | 2,54% | 41/1967 | 2,08% |
| | 3 | ND/1161 | <1% | 63/2226 | 2,83% | 33/1350 | 2,44% |
| **HA.143 0 A→G** | 1 | ND/1311 | <1% | ND/2870 | <1% | ND/1827 | <1% |
| | 2 | ND/1498 | <1% | 36/1924 | 1,87% | ND/1659 | <1% |
| | 3 | ND/955 | <1% | ND/2391 | <1% | ND/1452 | <1% |
| **HA.145 5 C→T** | 1 | ND/1333 | <1% | ND/2753 | <1% | 14/1233 | 1,14% |
| | 2 | ND/1846 | <1% | ND/2242 | <1% | ND/1895 | <1% |
| | 3 | ND/1093 | <1% | ND/2373 | <1% | ND/1449 | <1% |
| **HA.154 3 A→G** | 1 | 25/1209 | 2,07% | 94/2757 | 3,41% | 37/1142 | 3,24% |
| | 2 | ND/1660 | <1% | 56/1857 | 3,02% | 41/1585 | 2,59% |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | 3 | ND/1182 | <1% | ND/3324 | <1% | ND/1972 | <1% |
| **HA.162 4 C→A** | 1 | ND/998 | <1% | ND/2174 | <1% | ND/1478 | <1% |
| | 2 | ND/1173 | <1% | 25/1291 | 1,94% | ND/1120 | <1% |
| | 3 | ND/2218 | <1% | ND/3654 | <1% | ND/2244 | <1% |
| **HA.163 4 C→A** | 1 | ND/930 | <1% | ND/2032 | <1% | ND/1388 | <1% |
| | 2 | ND/1091 | <1% | 16/1218 | 1,31% | ND/1048 | <1% |
| | 3 | ND/2616 | <1% | ND/3704 | <1% | ND/2269 | <1% |
| **HA.163 8 C→A** | 1 | ND/932 | <1% | ND/1991 | <1% | ND/1368 | <1% |
| | 2 | ND/1083 | <1% | 15/1180 | 1,27% | ND/1010 | <1% |
| | 3 | ND/2600 | <1% | ND/3709 | <1% | ND/2276 | <1% |
| **HA.164 3 T→A** | 1 | ND/875 | <1% | ND/1892 | <1% | ND/1291 | <1% |
| | 2 | ND/1028 | <1% | 13/1110 | 1,17% | ND/944 | <1% |
| | 3 | ND/2612 | <1% | ND/3703 | <1% | ND/2278 | <1% |
| **HA.164 3 T→G** | 1 | ND/875 | <1% | ND/1892 | <1% | ND/1291 | <1% |
| | 2 | ND/1028 | <1% | 12/1110 | 1,08% | ND/944 | <1% |
| | 3 | ND/2612 | <1% | ND/3703 | <1% | ND/2278 | <1% |
| **HA.169 1 G→A** | 1 | ND/596 | <1% | ND/1110 | <1% | 7/404 | 1,73% |
| | 2 | ND/767 | <1% | ND/873 | <1% | ND/696 | <1% |
| | 3 | ND/2499 | <1% | ND/3575 | <1% | ND/2222 | <1% |
| **HA.169 3 A→T** | 1 | ND/582 | <1% | ND/1081 | <1% | 7/391 | 1,79% |
| | 2 | ND/751 | <1% | ND/864 | <1% | ND/690 | <1% |
| | 3 | ND/2310 | <1% | ND/3569 | <1% | ND/2219 | <1% |
| **HA.169 5 T→C** | 1 | ND/555 | <1% | ND/1030 | <1% | 7/366 | 1,91% |
| | 2 | ND/779 | <1% | ND/3557 | <1% | ND/688 | <1% |
| | 3 | ND/1767 | <1% | ND/3557 | <1% | ND/2220 | <1% |
| **HA.169 8 C→T** | 1 | ND/537 | <1% | ND/977 | <1% | ND/601 | <1% |
| | 2 | ND/758 | <1% | 11/852 | 1,29% | ND/681 | <1% |
| | 3 | ND/2260 | <1% | ND/3520 | <1% | ND/2113 | <1% |
| **HA.170 5 A→G** | 1 | ND/492 | <1% | ND/883 | <1% | ND/528 | <1% |
| | 2 | ND/733 | <1% | 11/832 | 1,32% | ND/660 | <1% |
| | 3 | ND/1709 | <1% | ND/3300 | <1% | ND/2016 | <1% |

531 Positions with a too low coverage (<298 reads/position) to detect ≥1% variants are marked

532 with an asterisk (*). Numbers are displayed as [number of variants]/[number of reads on that

533 position]. ND: not detected.

534

535 **Determining the influence of the minor variant detection method**

536    To isolate the effect of just the mSNV identification step in the DPP, independent of the

537    alignment step, quality-trimmed alignment files (*.bam files) of the data (subdivided per

538    virus, per SP and per DPP) were shared and subjected to the same DPP mSNV detection

539    process (in this case DPP3) and compared to the original outcomes from DPP1 and DPP2

540    (Table 5).  In the majority of positions, the different mSNV identification processes did not

541    influence the results, as 84% (119/142) of the mSNVs were identified regardless of the

542    mSNV identification process. Twenty-three mSNVs that were not reproduced by DPP3

543    mSNV identification analysis, were reproduced when the 'Direction and position Filters' in

544    DPP3 were ignored (Table 5, marked with # of ##). These parameters filter out mSNVs when

545    the set criteria for the read direction (variant must occur in both forward and reverse reads),

546    relative read direction (statistical approach of forward/reverse balance) and read position

547    (removal of systemic errors) are not met. However, DPP1 and DPP2 contain similar quality

548    parameters in their mSNV identification process, indicating that different DPPs deal

549    differently with quality parameters, and data could be excluded or included based on the DPP

550    used. In addition, 9 additional mSNVs were identified in the *.bam files compared to the

551    original mSNV outputs. It needs to be noted that the coverage of SP data analysed by DPP1

552    for positions identified with mSNVs was considerably lower compared to the coverage at that

553    position in the input *.bam files, suggesting additional quality filtering in the mSNV

554    detection step of DPP1. However, the influence on mSNV identification was limited most

555    likely due to the initial high nucleotide coverage.


556    To better visualise the differences in coverages and allele counts a graphical display of the

557    data for four positions showing mSNVs in different frequencies for each SP/DPP

558    combination is included in the supplemental material (S2 figure). In general, SNVs were

559    rarely missed due to low coverage, as also high coverage SP/DPP combinations display

560    discrepancies (table 3 and 4).

561

**Table 5. The reproducibility of positions with at least one ≥5% variant when alignment files from the respective DPPs are all uploaded into DPP3 for only the mSNV identification process versus when the mSNV identifications are fully performed by the respective DPPs.**

| Virus | Position | Sequence platform | Data Processing pipeline | | | | | | Bam file generating processing pipeline | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1 | | 2 | | 3 | | 1 | | 2 | | 3 | |
| | | | Minor variants | Percentage | Minor variants | Percentage | Minor variants | Percentage | Minor variants | Percentage | Minor variants | Percentage | Minor variants | Percentage |
| N L C H | PB2 .187 9 G→ A | 1 | 81/1301 | 6,2% | 246/2716 | 9,1% | 112/1203 | 9,3% | 132/1375 | 9,6% | 246/2716 | 9,1% | 121/1301 | 9,3% |
| | | 2 | 47/956 | 4,9% | 117/1137 | 10,3% | 114/1064 | 10,7% | 119/1122 | 10,6% | 117/1137 | 10,3% | 114/1064 | 10,7% |
| | | 3 | 49/530 | 9,2% | 131/1341 | 9,8% | 129/1338 | 9,6% | 54/542 | 10,0% | 131//1341 | 9,8% | 129/1338 | 9,6% |
| | PB2 .2101 G→ A | 1 | 53/1118 | 4,7% | 261/2704 | 9,7% | 110/897 | 12,3% | 138/1180 | 11,7% | 261/2704 | 9,7% | 121/1086 | 11,1% |
| | | 2 | 21/1578 | 1,3% | 125/1875 | 6,7% | 121/1463 | 8,3% | ND/1856## | <1% | ND/1850# | <1% | 121/1463 | 8,3% |
| | | 3 | 13/542 | 2,4% | 199/1433 | 13,9% | 199/1435 | 13,9% | 87/625 | 13,9% | 199/1433 | 13,9% | 199/1435 | 13,9% |
| | PB2 .2277 T→ G | 1 | ND/479 | <1% | 86/1008 | 8,5% | 33/190 | 17,4% | ND/849 | <1% | ND/1008## | <1% | 37/281 | 13,2% |
| | | 2 | ND/557 | <1% | ND/623 | <1% | ND/534 | <1% | ND/619 | <1% | ND/623 | <1% | ND/534 | <1% |
| | | 3 | ND/680 | <1% | ND/1117 | <1% | ND/1024 | <1% | ND/708 | <1% | ND/1117 | <1% | ND/1027 | <1% |
| | PB1 .87 A→ G | 1 | ND/818 | <1% | ND/1754 | <1% | ND/1114 | <1% | ND/1264 | <1% | ND/1753 | <1% | ND/1114 | <1% |
| | | 2 | 25/230 | 10,9% | ND/376 | <1% | ND/328 | <1% | ND/368## | <1% | ND/376 | <1% | ND/328 | <1% |
| | | 3 | ND/275 | <1% | ND/537 | <1% | ND/537 | <1% | ND/278 | <1% | ND/537 | <1% | ND/537 | <1% |
| | PB1 .2240 G→ C | 1 | ND/664 | <1% | 54/1341 | 4,0% | 38/418 | 9,1% | ND/1004 | <1% | ND/1341# | <1% | 46/486 | 9,5% |
| | | 2 | ND/1231 | <1% | ND/1271 | <1% | ND/1233 | <1% | ND/1277 | <1% | ND/1271 | <1% | ND/1235 | <1% |
| | | 3 | ND/161 | <1% | ND/277 | <1% | ND/276 | <1% | ND/163 | <1% | ND/277 | <1% | ND/276 | <1% |
| | PB1 .2268 A→ G | 1 | ND/336 | <1% | 29/641 | 4,5% | 11/176 | 6,3% | 15/322* | 4,66%* | 37/641 | 5,8% | 13/213 | 6,1% |
| | | 2 | ND/993 | <1% | ND/1026 | <1% | ND/1002 | <1% | ND/1025 | <1% | ND/1026 | <1% | ND/1002 | <1% |
| | | 3 | ND/53 | <1% | ND/159 | <1% | ND/148 | <1% | ND/90 | <1% | ND/159 | <1% | ND/151 | <1% |
| | PA. 2167 T→ G | 1 | ND/141 | <1% | ND/288 | <1% | ND/154 | <1% | ND/235 | <1% | 21/288* | 7,29%* | ND/154 | <1% |
| | | 2 | ND/757 | <1% | ND/807 | <1% | ND/773 | <1% | ND/812 | <1% | ND/807 | <1% | ND/773 | <1% |
| | | 3 | ND/704 | <1% | ND/1070 | <1% | ND/1077 | <1% | ND/714 | <1% | ND/1070 | <1% | ND/1078 | <1% |
| | HA. 104 A→ G | 1 | ND/733 | <1% | ND/1761 | <1% | ND/1151 | <1% | ND/1175 | <1% | ND/1761 | <1% | ND/1135 | <1% |
| | | 2 | ND/437 | <1% | ND/1370 | <1% | ND/1156 | <1% | ND/1326 | <1% | ND/1369 | <1% | ND/1142 | <1% |
| | | 3 | ND/1 | <1% | ND/1105 | <1% | 12/105 | 11,4% | ND/6 | <1% | ND/1105 | <1% | 12/105 | 11,4% |
| | HA. 168 | 1 | ND/390 | <1% | ND/694 | <1% | 11/217 | 5,1% | ND/610 | <1% | ND/694 | <1% | 13/260 | 5,0% |

| | Mutation | # | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **9 T→ C** | 2 | ND/2018 | <1% | ND/4083 | <1% | ND/3979 | <1% | ND/4045 | <1% | ND/4081 | <1% | ND/3979 | <1% |
| | | 3 | ND/937 | <1% | ND/1669 | <1% | ND/1680 | <1% | ND/1106 | <1% | ND/1669 | <1% | ND/1680 | <1% |
| | **NA.3 T→ C** | 1 | ND/32 | <1% | ND/105 | <1% | ND/49 | <1% | ND/92 | <1% | 7/105* | 6,67%* | ND/49 | <1% |
| | | 2 | ND/6 | <1% | ND/313 | <1% | ND/297 | <1% | ND/305 | <1% | ND/313 | <1% | ND/297 | <1% |
| | | 3 | ND/2 | <1% | ND/25 | <1% | ND/25 | <1% | ND/6 | <1% | ND/25 | <1% | ND/25 | <1% |
| | **NP.105 A→ G** | 1 | ND/182 | <1% | ND/449 | <1% | ND/343 | <1% | ND/374 | <1% | 6/449* | 1,34%* | ND/343 | <1% |
| | | 2 | 83/1507 | 5,5% | ND/1890 | <1% | ND/1804 | <1% | ND/1866## | <1% | ND/1890 | <1% | ND/1805 | <1% |
| | | 3 | ND/89 | <1% | ND/704 | <1% | ND/702 | <1% | ND/246 | <1% | ND/704 | <1% | ND/703 | <1% |
| | **NP.1239 A→ T** | 1 | 32/2428 | 1,3% | 279/5410 | 5,2% | ND/3092 | <1% | ND/3372## | <1% | ND/5410# | <1% | ND/3092 | <1% |
| | | 2 | ND/2345 | <1% | ND/2643 | <1% | ND/2453 | <1% | ND/2626 | <1% | ND/2643 | <1% | ND/2453 | <1% |
| | | 3 | ND/1711 | <1% | ND/2111 | <1% | ND/2117 | <1% | ND/1712 | <1% | ND/2111 | <1% | ND/2117 | <1% |
| | **NP.1489 G→ A** | 1 | ND/182 | <1% | 26/336 | 7,7% | ND/172 | <1% | ND/242 | <1% | 26/376* | 6,9% | ND/172 | <1% |
| | | 2 | ND/436 | <1% | ND/452 | <1% | ND/444 | <1% | ND/451 | <1% | ND/451 | <1% | ND/444 | <1% |
| | | 3 | ND/1320 | <1% | ND/1799 | <1% | ND/1799 | <1% | ND/1325 | <1% | ND/1799 | <1% | ND/1799 | <1% |
| | **NS.827 C→ T** | 1 | ND/249 | <1% | 19/419 | 4,5% | ND/205 | <1% | ND/365 | <1% | 21/412 | 5,3% | ND/205 | <1% |
| | | 2 | ND/1316 | <1% | ND/1423 | <1% | ND/1375 | <1% | ND/1427 | <1% | ND/1422 | <1% | ND/1375 | <1% |
| | | 3 | ND/2091 | <1% | ND/2901 | <1% | ND/2757 | <1% | ND/2293 | <1% | ND/2898 | <1% | ND/2929 | <1% |
| | **NS829 G→ T** | 1 | ND/221 | <1% | 19/380 | 5,0% | ND/179 | <1% | ND/328 | <1% | 19/376 | 5,4% | ND/179 | <1% |
| | | 2 | ND/1302 | <1% | ND/1391 | <1% | ND/1341 | <1% | ND/1388 | <1% | ND/1389 | <1% | ND/1341 | <1% |
| | | 3 | ND/2117 | <1% | ND/2852 | <1% | ND/2727 | <1% | ND/2279 | <1% | ND/2852 | <1% | ND/2880 | <1% |
| | **NS.833 A→ T** | 1 | ND/187 | <1% | ND/287 | <1% | 5/88 | 5,7% | ND/259 | <1% | 11/257* | 4,28%* | 5/96 | 5,2% |
| | | 2 | ND/1224 | <1% | ND/1327 | <1% | ND/1284 | <1% | ND/1314 | <1% | ND/1322 | <1% | ND/1284 | <1% |
| | | 3 | ND/1367 | <1% | ND/2430 | <1% | ND/2333 | <1% | ND/1779 | <1% | ND/2430 | <1% | ND/2360 | <1% |
| **D ET U** | **PB2.900 A→ G** | 1 | 38/1335 | 2,9% | 136/2740 | 5,0% | 61/1231 | 5,0% | 68/1328 | 5,12 | 136/2740 | 4,96 | 65/1322 | 4,92 |
| | | 2 | 35/1645 | 2,1% | 77/1800 | 4,3% | 66/1629 | 4,1% | 70/1775 | 4,0% | 77/1800 | 4,3% | 66/1629 | 4,1% |
| | | 3 | 30/861 | 3,5% | 86/2308 | 3,7% | 47/1245 | 3,8% | ND/1001## | <1% | ND/2308# | <1% | 47/1245 | 3,8% |
| | **PB2.1054 T→ C** | 1 | 69/1369 | 5,0% | 168/2637 | 6,4% | 97/1304 | 7,4% | 105/1393 | 7,5% | 168/2637 | 6,4% | 100/1376 | 7,3% |
| | | 2 | 60/1477 | 4,1% | 115/1836 | 6,3% | 99/1605 | 6,2% | 113/1810 | 6,2% | 115/1836 | 6,3% | 99/1605 | 6,2% |
| | | 3 | 6/392 | 1,5% | 94/2038 | 4,6% | 48/1054 | 4,6% | 32/524 | 6,1% | 94/2038 | 4,6% | 48/1054 | 4,6% |
| | **PB2.2257 A→ C** | 1 | ND/867 | <1% | ND/1563 | <1% | 24/463 | 5,2% | ND/1447 | <1% | ND/1562 | <1% | 26/472 | 5,5% |
| | | 2 | ND/531 | <1% | ND/581 | <1% | ND/378 | <1% | ND/588 | <1% | ND/580 | <1% | ND/378 | <1% |
| | | 3 | ND/893 | <1% | ND/2286 | <1% | ND/1346 | <1% | ND/1341 | <1% | ND/2185 | <1% | ND/1347 | <1% |
| | **PB2.2277 T→ G** | 1 | ND/644 | <1% | 52/1150 | 4,5% | 27/307 | 8,8% | ND/1062 | <1% | ND/1150# | <1% | 28/381 | 7,4% |
| | | 2 | ND/418 | <1% | ND/472 | <1% | ND/284 | <1% | ND/474 | <1% | ND/472 | <1% | ND/284 | <1% |
| | | 3 | ND/1208 | <1% | ND/1948 | <1% | ND/1209 | <1% | ND/1251 | <1% | ND/1948 | <1% | ND/1214 | <1% |
| | **PB1.14 C→ T** | 1 | ND/144 | <1% | 48/433 | 11,1% | ND/239 | <1% | ND/362 | <1% | 48/433 | 11,1% | ND/239 | <1% |
| | | 2 | ND/90 | <1% | ND/355 | <1% | ND/304 | <1% | ND/345 | <1% | ND/351 | <1% | ND/304 | <1% |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 3 | ND/562 | <1% | ND/792 | <1% | ND/496 | <1% | ND/633 | <1% | ND/655 | <1% | ND/504 | <1% |
| **PB1.23 T→G** | 1 | ND/207 | <1% | 30/535 | 5,6% | ND/315 | <1% | ND/470 | <1% | 30/535 | 5,6% | ND/315 | <1% |
| | 2 | ND/103 | <1% | ND/365 | <1% | ND/319 | <1% | ND/365 | <1% | 4/365* | 1,96%* | ND/319 | <1% |
| | 3 | ND/699 | <1% | ND/950 | <1% | ND/609 | <1% | ND/702 | <1% | ND/950 | <1% | ND/609 | <1% |
| **PB1.87 A→G** | 1 | ND/744 | <1% | ND/1644 | <1% | ND/1076 | <1% | ND/1218 | <1% | ND/1644 | <1% | ND/1076 | <1% |
| | 2 | 49/365 | 13,4% | ND/677 | <1% | ND/576 | <1% | 13/638 | 2,0% | ND/674 | <1% | ND/576 | <1% |
| | 3 | ND/721 | <1% | ND/1156 | <1% | ND/793 | <1% | ND/731 | <1% | ND/1156 | <1% | ND/793 | <1% |
| **PB1.2240 G→C** | 1 | ND/757 | <1% | 23/1517 | 1,5% | 26/515 | 5,0% | ND/1266 | <1% | ND/1515# | <1% | 28/631 | 4,4% |
| | 2 | ND/944 | <1% | ND/985 | <1% | ND/806 | <1% | ND/994 | <1% | ND/984 | <1% | ND/806 | <1% |
| | 3 | ND/274 | <1% | ND/439 | <1% | ND/253 | <1% | ND/301 | <1% | ND/439 | <1% | ND/253 | <1% |
| **PB1.2268 A→G** | 1 | 5/470 | 1,1% | 33/928 | 3,6% | 22/278 | 7,9% | 28/420 | 6,7% | ND/928## | <1% | 23/354 | 6,5% |
| | 2 | ND/798 | <1% | ND/829 | <1% | ND/671 | <1% | ND/839 | <1% | ND/829 | <1% | ND/671 | <1% |
| | 3 | ND/109 | <1% | ND/259 | <1% | ND/123 | <1% | ND/193 | <1% | ND/259 | <1% | ND/126 | <1% |
| **PB1.2271 A→G** | 1 | 12/446 | 2,7% | 59/901 | 6,5% | 16/263 | 6,1% | 29/413 | 7,0% | 59/901 | 6,6% | 21/336 | 6,3% |
| | 2 | ND/729 | <1% | 47/810 | 5,8% | 40/649 | 6,2% | 43/750* | 5,73%* | 47/810 | 5,8% | 40/649 | 6,2% |
| | 3 | 1/32 | 3,1% | ND/123 | <1% | 2/83 | 2,4% | 5/75 | 6,7% | 5/124* | 4,03%* | 2/83 | 2,4% |
| **HA.867 C→T** | 1 | 59/1533 | 3,8% | 206/3183 | 6,5% | 104/1537 | 6,8% | 112/1584 | 7,1% | 206/3183 | 6,5% | 109/1573 | 6,9% |
| | 2 | 59/2031 | 2,9% | 150/2525 | 5,9% | 127/2253 | 5,6% | 144/2502 | 5,8% | 150/2525 | 5,9% | 127/2253 | 5,6% |
| | 3 | 11/180 | 6,1% | 48/647 | 7,4% | 28/385 | 7,3% | 13/182 | 7,1% | 48/647 | 7,4% | 28/385 | 7,3% |
| **HA.963 T→C** | 1 | 122/1401 | 8,7% | 446/3071 | 14,5% | 189/1419 | 13,3% | 200/1468 | 13,6% | 446/3071 | 14,5% | 193/1455 | 13,3% |
| | 2 | 90/1517 | 5,9% | 318/2189 | 14,5% | 247/1828 | 13,5% | 308/2165 | 14,2% | 318/2189 | 14,5% | 247/1828 | 13,5% |
| | 3 | 5/69 | 7,2% | 107/606 | 17,7% | 47/293 | 16,0% | 12/81 | 14,8% | 107/606 | 17,7% | 47/293 | 16,0% |
| **NP.1491 C→A** | 1 | ND/278 | <1% | 71/583 | 12,2% | ND/206 | <1% | ND/390 | <1% | ND/579# | <1% | ND/206 | <1% |
| | 2 | ND/723 | <1% | ND/769 | <1% | ND/692 | <1% | ND/766 | <1% | ND/769 | <1% | ND/692 | <1% |
| | 3 | ND/799 | <1% | ND/2031 | <1% | ND/1206 | <1% | ND/858 | <1% | ND/2031 | <1% | ND/1206 | <1% |
| **NA.65 T→C** | 1 | 19/503 | 3,8% | 52/1229 | 4,2% | 16/467 | 3,4% | 22/535 | 4,1% | 52/1229 | 4,2% | 20/540 | 3,7% |
| | 2 | 20/662 | 3,0% | 50/1104 | 4,5% | 45/992 | 4,5% | 52/1063 | 4,9% | 50/1104 | 4,5% | 45/992 | 4,5% |
| | 3 | 24/557 | 4,3% | 53/1099 | 4,8% | 37/727 | 5,1% | 28/584 | 4,8% | 53/1099 | 4,8% | 37/727 | 5,1% |
| **NA.78 T→C** | 1 | 23/599 | 3,8% | 57/1403 | 4,1% | 20/557 | 3,6% | 23/622 | 3,7% | 57/1403 | 4,1% | 24/638 | 3,8% |
| | 2 | 21/692 | 3,0% | 55/1147 | 4,8% | 50/1033 | 4,8% | 54/1109 | 4,9% | 55/1147 | 4,8% | 50/1033 | 4,8% |
| | 3 | 23/580 | 4,0% | 51/1124 | 4,5% | 37/735 | 5,0% | 27/585 | 4,6% | ND/1124# | <1% | 37/735 | 5,0% |
| **NA.89 T→C** | 1 | 23/713 | 3,2% | 55/1670 | 3,3% | 22/651 | 3,4% | 26/731 | 3,6% | 55/1670 | 3,3% | 26/751 | 3,5% |
| | 2 | 23/798 | 2,9% | 56/1261 | 4,4% | 50/1134 | 4,4% | 54/1224 | 4,4% | 56/1261 | 4,4% | 50/1134 | 4,4% |
| | 3 | 24/580 | 4,1% | 55/1196 | 4,6% | 40/775 | 5,2% | 28/587 | 4,8% | 55/1196 | 4,6% | 40/775 | 5,2% |
| **NA.117 T→C** | 1 | 37/908 | 4,1% | 87/2140 | 4,1% | 36/818 | 4,4% | 40/914 | 4,4% | 87/2140 | 4,7% | 43/922 | 4,7% |
| | 2 | 28/1102 | 2,5% | 67/1631 | 4,1% | ND/1459 | <1% | 70/1586 | 4,4% | 67/1631 | 4,1% | ND/1459 | <1% |

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 3 | 22/531 | 4,1% | 57/1276 | 4,5% | 42/812 | 5,2% | 28/544 | 5,2% | ND/1276# | <1% | 42/812 | 5,2% |
| | **NA. 126 T→ C** | 1 | 37/983 | 3,8% | 83/2294 | 3,6% | 36/876 | 4,1% | 39/973 | 4,0% | 83/2294 | 3,6% | 43/981 | 4,4% |
| | | 2 | 31/1126 | 2,8% | 72/1676 | 4,3% | 65/1502 | 4,3% | 75/1616 | 4,6% | 72/1676 | 4,3% | 65/1502 | 4,3% |
| | | 3 | 26/519 | 5,0% | 62/1395 | 4,4% | 43/812 | 5,3% | 30/537 | 5,6% | 62/1395 | 4,4% | 43/812 | 5,3% |
| **U K D D** | **PB2 .227 7 T→ G** | 1 | ND/415 | <1% | 28/507 | 5,5% | ND/475 | <1% | ND/503 | <1% | ND/507# | <1% | ND/475 | <1% |
| | | 2 | ND/589 | <1% | ND/620 | <1% | ND/601 | <1% | ND/627 | <1% | ND/620 | <1% | ND/601 | <1% |
| | | 3 | ND/1140 | <1% | ND/1996 | <1% | ND/2065 | <1% | ND/1186 | <1% | ND/1996 | <1% | ND/2071 | <1% |
| | **PB2 .227 8 T→ G** | 1 | ND/367 | <1% | ND/471 | <1% | ND/464## | <1% | ND/465 | <1% | ND/471 | <1% | 17/268 | 6,3% |
| | | 2 | ND/581 | <1% | ND/613 | <1% | ND/581 | <1% | ND/621 | <1% | ND/588 | <1% | ND/581 | <1% |
| | | 3 | ND/1141 | <1% | ND/1985 | <1% | ND/1993 | <1% | ND/1184 | <1% | ND/1975 | <1% | ND/2004 | <1% |
| | **PB1 .87 A→ G** | 1 | ND/387 | <1% | ND/440 | <1% | ND/439 | <1% | ND/451 | <1% | ND/417 | <1% | ND/439 | <1% |
| | | 2 | 26/327 | 8,0% | 32/395 | 8,1% | ND/351 | <1% | 33/385 | 8,6% | ND/395# | <1% | ND/351 | <1% |
| | | 3 | ND/617 | <1% | ND/1133 | <1% | ND/1136 | <1% | ND/622 | <1% | ND/1133 | <1% | ND/1136 | <1% |
| | **PB1 .728 C→ A** | 1 | ND/750 | <1% | ND/832 | <1% | ND/836 | <1% | ND/853 | <1% | ND/832 | <1% | ND/836 | <1% |
| | | 2 | ND/776 | <1% | 52/928 | 5,6% | ND/829 | <1% | ND/888 | <1% | ND/912## | <1% | ND/829 | <1% |
| | | 3 | ND/2459 | <1% | ND/4290 | <1% | ND/4293 | <1% | ND/2471 | <1% | ND/4287 | <1% | ND/4292 | <1% |
| | **PB1 .730 C→ T** | 1 | ND/742 | <1% | ND/824 | <1% | ND/826 | <1% | ND/844 | <1% | ND/824 | <1% | ND/826 | <1% |
| | | 2 | ND/767 | <1% | 57/1008 | 5,7% | ND/832 | <1% | ND/893 | <1% | ND/1008# | <1% | ND/832 | <1% |
| | | 3 | ND/2339 | <1% | ND//4286 | <1% | ND/4289 | <1% | ND/2464 | <1% | ND/4285 | <1% | ND/4284 | <1% |
| | **PB1 .883 G→ C** | 1 | ND/942 | <1% | ND/997 | <1% | ND/997 | <1% | ND/1016 | <1% | ND/997 | <1% | ND/997 | <1% |
| | | 2 | ND/1689 | <1% | ND/1856 | <1% | ND/1760 | <1% | ND/1867 | <1% | ND/1856 | <1% | ND/1760 | <1% |
| | | 3 | ND/2479 | <1% | 47/690 | 6,8% | ND/3681 | <1% | ND/2635 | <1% | ND/690## | <1% | ND/3697 | <1% |
| | **PA. 49 G→ C** | 1 | ND/103 | <1% | 6/117 | 5,1% | ND/115 | <1% | ND/113 | <1% | ND/117# | <1% | ND/115 | <1% |
| | | 2 | ND/337 | <1% | ND/435 | <1% | ND/392 | <1% | ND/441 | <1% | ND/434 | <1% | ND/392 | <1% |
| | | 3 | ND/111 | <1% | ND/207 | <1% | ND/204 | <1% | ND/113 | <1% | ND/206 | <1% | ND/206 | <1% |
| | **PA. 82 C→ T** | 1 | ND/155 | <1% | ND/180 | <1% | ND/177 | <1% | ND/179 | <1% | ND/180 | <1% | ND/177 | <1% |
| | | 2 | ND/695 | <1% | ND/809 | <1% | ND/745 | <1% | ND/797 | <1% | ND/809 | <1% | ND/745 | <1% |
| | | 3 | ND/64 | <1% | ND/247 | <1% | 30/248 | 12,1% | ND/74 | <1% | ND/247 | <1% | 30/248 | 12,1% |
| | **NS. 811 G→ T** | 1 | ND/221 | <1% | 17/270 | 6,3% | ND/249 | <1% | ND/261 | <1% | ND/270# | <1% | ND/249 | <1% |
| | | 2 | ND/2452 | <1% | ND/2725 | <1% | ND/2557 | <1% | ND/2742 | <1% | ND/2725 | <1% | ND/2557 | <1% |
| | | 3 | ND/3117 | <1% | ND/4125 | <1% | ND/4139 | <1% | ND/3188 | <1% | ND/4124 | <1% | ND/4142 | <1% |

*Locations containing mSNV detections in the DPP3 mSNV analysis of the bam files but not

in the original DPPs; Locations containing ≥1% mSNVs that could be reproduced by deleting

DPP3s default 'Direction and position filters' with those exactly reproduced (#) and those

approximately reproduced but with different coverages and/or variants (##).

570

# **Discussion**

NGS data are used for different applications. Although sequence technologies and the

accompanying analysis tools are subjected to rapid development, a lot of follow-up research

is based on initial findings. Accuracy and repeatability are key values for proper scientific

research but the impact of NGS results also reaches beyond science to clinical settings where

important clinical management and treatment decisions are based on such results. In this

study the comparability of NGS data analyses were analysed using identical input material

per virus but different laboratory workflows from nucleic acid extraction and sequencing to

data analysis. In addition, the COMPARE "Data Hub" platform was tested for the purpose of

sharing large raw datafiles between institutions in an outbreak situation. Using this platform,

raw sequence data files up to the size of 8 Gigabytes, alignment files and metadata files of

three influenza A/H5N8 viruses were successfully shared in real-time among 3 institutions to

allow independent sequencing and analysis procedures, including mSNV identification, to be

performed. The Data Hub is available to all institutions.

The aim of this study was to determine how comparable consensus and minority variant

results were between laboratories performing their standard analyses, and whether

discrepancies could be attributed to the SP, DPP or a combination of both. With the lack of a

ground truth/gold standard, all data obtained were compared amongst each other.

Importantly, reliable consensus sequences were generated independently of the SP/DPP

combination used, although the well-known artefactual InDels in homopolymer regions in

SP3 (Roche 454 genome sequencer) sequence data required manual editing. Such consensus

sequences routinely form the basis for a detailed characterization of the influenza strain in an

593   outbreak situation, as they are used for the prediction of pathogenicity and pandemic potential

594   of influenza strains.

595   In contrast to the reproducible generation of consensus genome sequences, the hypothesis

596   that minority variants could be identified reproducibly has to be rejected. The observed

597   differences were mainly attributed to the alignment processes in the different DPPs. The

598   interpretation of minority variant analysis thus needs a different level of careful

599   standardization and awareness about the possible limitations as shown in this study.

600   Reproducibility of mSNV results appeared to be influenced by both the different SPs

601   (resulting in different sequence depths Fig. 2) and DPPs (resulting in differences in alignment

602   and mSNV identification of the same input data, Fig. 2 and Table 5) . There was limited

603   reproducibility of mSNV identification data, even for relative high frequency mSNVs. As

604   expected, the reproducibility was best (30%) for mSNVs occurring in high frequency

605   ($\geq$10%), and least for the low frequent ($\geq$1%) mSNVs (9.4% to 31.1%). Also, the number of

606   positions with 1-5% mSNVs (with sufficient coverage) was much higher (250 in SP1 data,

607   213 in SP2 data, and 45 in SP3 data) than the number of positions with >5-10% mSNVs

608   (n=27) or >10% mSNVs (n=10).

609   The set-up of this study allowed many variables to influence the final result. The differences

610   from first laboratory procedures and sample preparations up to the final analysis methods can

611   all have contributed to the observed differences in mSNV identification. At this level,

612   especially with lacking an NGS gold standard, it becomes difficult to determine which

613   identified mSNVs are 'true variants' and which could be due to systematic errors introduced

614   by RNA isolation methods, amplification, sequencing or manipulated by data processing

615   pipeline settings. Unsurprisingly, the results of this study imply that the choice of SP

616   influences the final output, but the results from this study also indicate that the DPP,

617 especially the alignment process, influences coverage. The SP and DPP derived differences in

618 coverage are of importance because up to a certain (currently unknown, probably SP/DPP

619 dependent) threshold, a higher coverage will provide a more reliable result about the presence

620 of mSNVs. Although the aim of this study was to explicitly compare the three institutions

621 own standard workflows, some parameters (like the phred score and detection limit) were

622 synchronized between the different DPPs. Moreover, the data from each SP were re-

623 processed in each DPP. However, all DPPs use different underlying algorithms and interpret

624 the set parameters differently which might all contribute to the observed differences. These

625 results are partly in line with previous research that showed the need of NGS result validation

626 and concluded that only those mSNVs with a coverage >100 and a frequency of >40% could

627 be identified by NGS methods without secondary confirmation [32], however, this conclusion

628 was based on using the same sample preparation method within a single laboratory. Another

629 recent study sets the cut-off for intrahost virus diversity at 3% with input of at least 1000

630 RNA copies and a read depth of at least 400x at each genome position for Illumina

631 sequencing [33].

632 Although some studies have been published on SP error rates [34-37] and PCR amplification

633 induced variants [38-41], a gold standard system for mSNV analysis is lacking. In addition,

634 the DPPs can alter the data due to elimination or inclusion of certain sequences based on the

635 set quality parameters. Allowing too many low-quality reads or being too stringent on the

636 data will influence the coverage per position and might also influence the accuracy of the

637 mSNV identification rate, especially when the coverage is low [42, 43]. Although a low

638 comparability of mSNVs identified in the different SP and DPP combinations was observed,

639 it can be concluded that 454 (SP3) sequencing has approximately the same accuracy as

640 Illumina (SP1 and 2) sequencing based on the number and percentage of reproducible

641 mSNVs in this dataset when ignoring InDel errors in homopolymer regions. Although, Roche

642     454 sequencing machines are no longer in production, it added value to include 454

643     sequencing as an alternative sequence platform with alternative chemistry to Illumina. In

644     addition, because Roche 454 was the first commercially successful next generation

645     sequencing system, it was used in research that served as a fundament for follow-up studies

646     [44]. A comparison of Illumina with newer third or fourth generation sequencing platforms

647     (e.g. Nanopore or Pac Bio) would be interesting in the future. However, the overall error rate

648     remains higher than the shorter read technologies and recent work concludes that these new

649     platforms are currently not suitable for the detection of minor variants [33]. In addition, it

650     would be interesting to compare mSNV results of SPs outputting small sequence reads (like

651     Illumina, 454 and Ion Torrent) to new sequencing techniques that output full-length sequence

652     data (e.g. Nanopore [45]). The latter might be less vulnerable to quality trimming parameters

653     compared to small reads and might provide a more consistent nucleotide coverage over

654     complete gene segment.

655     For mSNV analyses by different labs, very stringent SP/DPP protocols need to be evaluated,

656     for instance by cross-validating results. To allow a better comparison it would be

657     recommended to create some kind of gold standard by for instance evaluating parameters

658     based on sequencing of technical replicates, and controlled mixes of clones. The mSNV

659     analysis can be valuable for epidemiological tracing, to monitor early evolutionary events, or

660     drug resistance, possibly host adaptation, but this would require reproducibility of study

661     outcomes within and between laboratories. As this is currently not that case, more

662     understanding of biases and errors generated by sample processing (enrichment procedures),

663     sequencing strategy (amplicons, shotgun), sequencing chemistry (each of which have their

664     own internal error rates) and the approach to data processing and analysis is needed.

665     Understanding the parameters and thresholds in the software can be difficult and a systematic

666     study using a pipeline where the effect of changing each of these parameters both

667 individually and in combination is required to determine the optimal settings for minor

668 variant analysis.

669 As alternate high-throughput sequencing technologies arise there will be a need to understand

670 inherent error profiles and how those are handled in data processing approaches. Cross-

671 validation should be supported by international proficiency tests on NGS techniques

672 including mSNV analyses that would be instrumental in validation of results and may foster

673 the trust in NGS-based diagnostics.

# Acknowledgements

# References

678 1. Heather, J.M. and B. Chain, *The sequence of sequencers: The history of sequencing DNA.*
679    Genomics, 2016. **107**(1): p. 1-8.
680 2. van Dijk, E.L., et al., *Ten years of next-generation sequencing technology.* Trends Genet,
681    2014. **30**(9): p. 418-26.
682 3. Ekblom, R. and J. Galindo, *Applications of next generation sequencing in molecular ecology of*
683    *non-model organisms.* Heredity (Edinb), 2011. **107**(1): p. 1-15.
684 4. Koser, C.U., et al., *Rapid whole-genome sequencing for investigation of a neonatal MRSA*
685    *outbreak.* N Engl J Med, 2012. **366**(24): p. 2267-75.
686 5. Mellmann, A., et al., *Prospective genomic characterization of the German enterohemorrhagic*
687    *Escherichia coli O104:H4 outbreak by rapid next generation sequencing technology.* PLoS
688    One, 2011. **6**(7): p. e22751.
689 6. Leitner, T., et al., *Analysis of heterogeneous viral populations by direct DNA sequencing.*
690    Biotechniques, 1993. **15**(1): p. 120-7.
691 7. Tsiatis, A.C., et al., *Comparison of Sanger sequencing, pyrosequencing, and melting curve*
692    *analysis for the detection of KRAS mutations: diagnostic and clinical implications.* J Mol
693    Diagn, 2010. **12**(4): p. 425-32.
694 8. Glenn, T.C., *Field guide to next-generation DNA sequencers.* Mol Ecol Resour, 2011. **11**(5): p.
695    759-69.
696 9. Li, Y., et al., *Ion Torrent™ Next Generation Sequencing – Detect 0.1% Low Frequency Somatic*
697    *Variants and Copy Number Variations simultaneously in Cell-Free DNA.* Thermo Fisher
698    Scientific, 2017.

699  10.  Schirmer, M., et al., *Illumina error profiles: resolving fine-scale variation in metagenomic*
700      *sequencing data.* BMC Bioinformatics, 2016. **17**: p. 125.
701  11.  Lou, D.I., et al., *High-throughput DNA sequencing errors are reduced by orders of magnitude*
702      *using circle sequencing.* Proc Natl Acad Sci U S A, 2013. **110**(49): p. 19872-7.
703  12.  World Organisation for Animal Health, O.I.E., *Update on highly pathogenic avian influenza in*
704      *animals (typeH5 and H7).* 2014.
705  13.  World Organisation for Animal Health, O.I.E., *Update on highly pathogenic avian influenza in*
706      *animals (typeH5 and H7).* 2015.
707  14.  Hanna, A., et al., *Genetic Characterization of Highly Pathogenic Avian Influenza (H5N8) Virus*
708      *from Domestic Ducks, England, November 2014.* Emerg Infect Dis, 2015. **21**(5): p. 879-82.
709  15.  Harder, T., et al., *Influenza A(H5N8) Virus Similar to Strain in Korea Causing Highly*
710      *Pathogenic Avian Influenza in Germany.* Emerg Infect Dis, 2015. **21**(5): p. 860-3.
711  16.  Bouwstra, R., et al., *Full-Genome Sequence of Influenza A(H5N8) Virus in Poultry Linked to*
712      *Sequences of Strains from Asia, the Netherlands, 2014.* Emerg Infect Dis, 2015. **21**(5): p. 872-
713      4.
714  17.  Verhagen, J.H., et al., *Wild bird surveillance around outbreaks of highly pathogenic avian*
715      *influenza A(H5N8) virus in the Netherlands, 2014, within the context of global flyways.* Euro
716      Surveill, 2015. **20**(12).
717  18.  Poen, M.J., et al., *Local amplification of highly pathogenic avian influenza H5N8 viruses in*
718      *wild birds in the Netherlands, 2016 to 2017.* Euro Surveill, 2018. **23**(4).
719  19.  Global Consortium for, H.N. and V. Related Influenza, *Role for migratory wild birds in the*
720      *global spread of avian influenza H5N8.* Science, 2016. **354**(6309): p. 213-217.
721  20.  Harrison, P.W., et al., *The European Nucleotide Archive in 2018.* Nucleic Acids Research,
722      2019. **47**(D1): p. D84-D88.
723  21.  Karsch-Mizrachi, I., et al., *The international nucleotide sequence database collaboration.*
724      Nucleic Acids Research, 2018. **46**(D1): p. D48-D51.
725  22.  Amid, C., et al., *The COMPARE Data Hubs.* bioRxiv, 2019: p. 555938.
726  23.  Richard, M., et al., *Low Virulence and Lack of Airborne Transmission of the Dutch Highly*
727      *Pathogenic Avian Influenza Virus H5N8 in Ferrets.* PLoS One, 2015. **10**(6): p. e0129827.
728  24.  Linster, M., et al., *Identification, characterization, and natural selection of mutations driving*
729      *airborne transmission of A/H5N1 virus.* Cell, 2014. **157**(2): p. 329-339.
730  25.  Li, H., *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.* arXiv,
731      2013.
732  26.  Li, H., et al., *The Sequence Alignment/Map format and SAMtools.* Bioinformatics, 2009.
733      **25**(16): p. 2078-9.
734  27.  Zerbino, D.R. and E. Birney, *Velvet: algorithms for de novo short read assembly using de*
735      *Bruijn graphs.* Genome Res, 2008. **18**(5): p. 821-9.
736  28.  Camacho, C., et al., *BLAST+: architecture and applications.* BMC Bioinformatics, 2009. **10**: p.
737      421.
738  29.  Hwang, S., et al., *Systematic comparison of variant calling pipelines using gold standard*
739      *personal exome variants.* Sci Rep, 2015. **5**: p. 17875.
740  30.  Hall, T.A., *BioEdit: a user-friendly biological sequence alignment editor and analysis program*
741      *for Windows 95/98/NT.* Nucleic Acids Symposium Series, 1999. **41**: p. 95-98.
742  31.  Dell, R.B., S. Holleran, and R. Ramakrishnan, *Sample size determination.* ILAR J, 2002. **43**(4):
743      p. 207-13.
744  32.  Mu, W., et al., *Sanger Confirmation Is Required to Achieve Optimal Sensitivity and Specificity*
745      *in Next-Generation Sequencing Panel Testing.* J Mol Diagn, 2016. **18**(6): p. 923-932.
746  33.  Grubaugh, N.D., et al., *An amplicon-based sequencing framework for accurately measuring*
747      *intrahost virus diversity using PrimalSeq and iVar.* Genome Biol, 2019. **20**(1): p. 8.
748  34.  Golan, D. and P. Medvedev, *Using state machines to model the Ion Torrent sequencing*
749      *process and to improve read error rates.* Bioinformatics, 2013. **29**(13): p. i344-51.

750  35.  Manley, L.J., D. Ma, and S.S. Levine, *Monitoring Error Rates In Illumina Sequencing.* J Biomol
751       Tech, 2016. **27**(4): p. 125-128.
752  36.  Nakamura, K., et al., *Sequence-specific error profile of Illumina sequencers.* Nucleic Acids Res,
753       2011. **39**(13): p. e90.
754  37.  Shao, W., et al., *Analysis of 454 sequencing error rate, error sources, and artifact
755       recombination for detection of Low-frequency drug resistance mutations in HIV-1 DNA.*
756       Retrovirology, 2013. **10**: p. 18.
757  38.  Acinas, S.G., et al., *PCR-induced sequence artifacts and bias: insights from comparison of two
758       16S rRNA clone libraries constructed from the same sample.* Appl Environ Microbiol, 2005.
759       **71**(12): p. 8966-9.
760  39.  Gorzer, I., et al., *The impact of PCR-generated recombination on diversity estimation of
761       mixed viral populations by deep sequencing.* J Virol Methods, 2010. **169**(1): p. 248-52.
762  40.  Judo, M.S., A.B. Wedel, and C. Wilson, *Stimulation and suppression of PCR-mediated
763       recombination.* Nucleic Acids Res, 1998. **26**(7): p. 1819-25.
764  41.  Meyerhans, A., J.P. Vartanian, and S. Wain-Hobson, *DNA recombination during PCR.* Nucleic
765       Acids Res, 1990. **18**(7): p. 1687-91.
766  42.  Quail, M.A., et al., *A tale of three next generation sequencing platforms: comparison of Ion
767       Torrent, Pacific Biosciences and Illumina MiSeq sequencers.* BMC Genomics, 2012. **13**: p. 341.
768  43.  Sims, D., et al., *Sequencing depth and coverage: key considerations in genomic analyses.* Nat
769       Rev Genet, 2014. **15**(2): p. 121-32.
770  44.  Liu, L., et al., *Comparison of next-generation sequencing systems.* J Biomed Biotechnol, 2012.
771       **2012**: p. 251364.
772  45.  Keller, M.W., et al., *Direct RNA Sequencing of the Coding Complete Influenza A Virus
773       Genome.* Sci Rep, 2018. **8**(1): p. 14408.

774

# Supporting information

776  **S1 Table. PCR primers used in SP3 to cover the influenza A H5N8 gene segments**

777
778  **S2 Table. SP/DPP overarching consensus sequences**

779
780  **S3 Table. Number of raw sequences and influenza virus reads per SP per virus**

781  **S1 File. DPP3 Sequence analysis protocol**

782
783  **S1 Figure. Nucleotide coverage.** The non-normalised nucleotide coverage displayed as

784  number of nucleotides per position for full genome sequences of the UKDD and DETU virus

785  reads mapped to the corresponding reference sequences. Panel A shows the coverage results

786  for the same SP dataset in the three different DPPs (DPP1: purple; DPP2: orange; DPP3 grey)

787  for each of the SP datasets. Panel B shows the coverage when the same DPP is used to

788  analyse data from the three different SPs (SP1: lilac; SP2: yellow; SP3:green) for each of the

789    DPPs. The X-axis represents the position in the genome, the Y-axis represents the number of

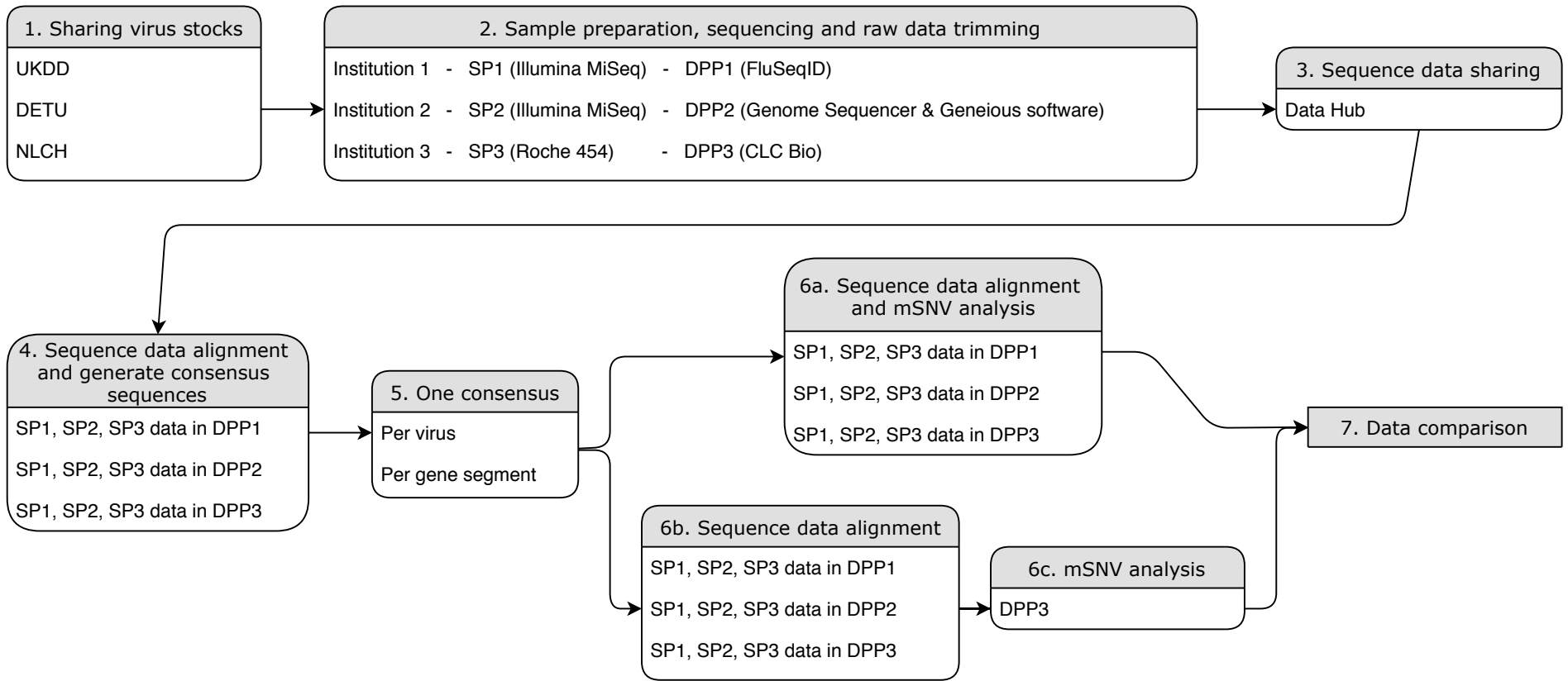790    sequence reads per position.

791

792    **S2 Figure. Graphical display of the coverage and allele counts for four positions,**

793    **showing mSNVs in different frequencies for each SP/DPP combination.** Arrows indicate

794    the approximate percentages in which the mSNVs were detected; 1-5% (orange), 5-10%
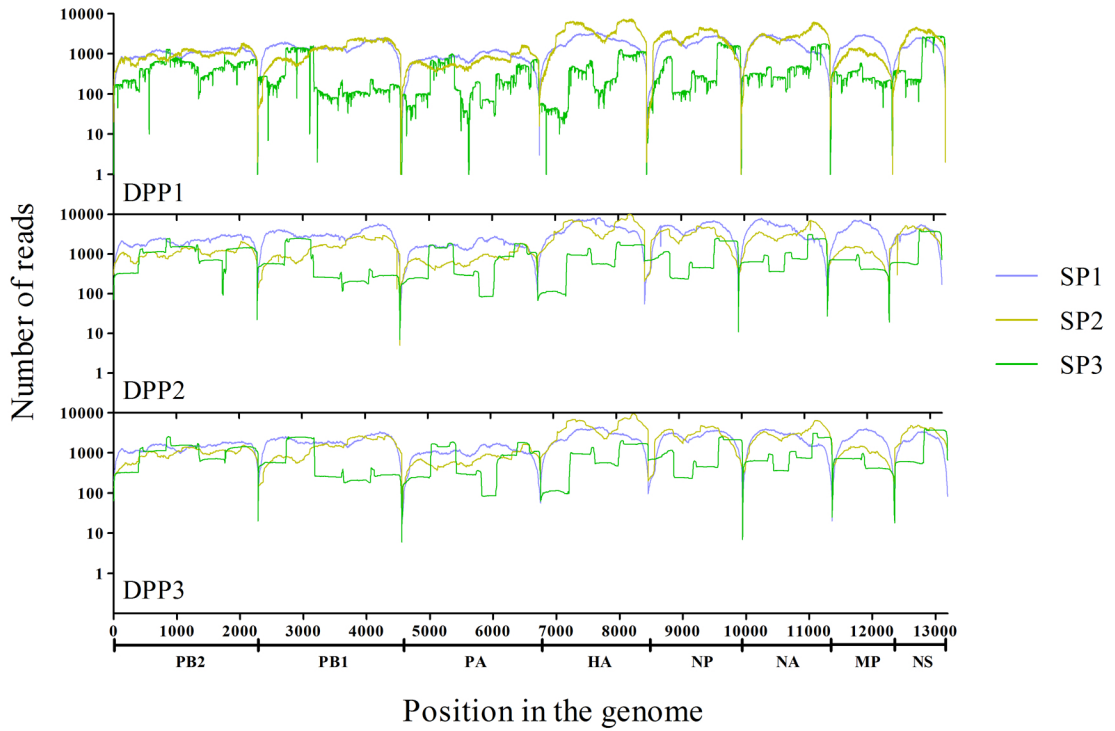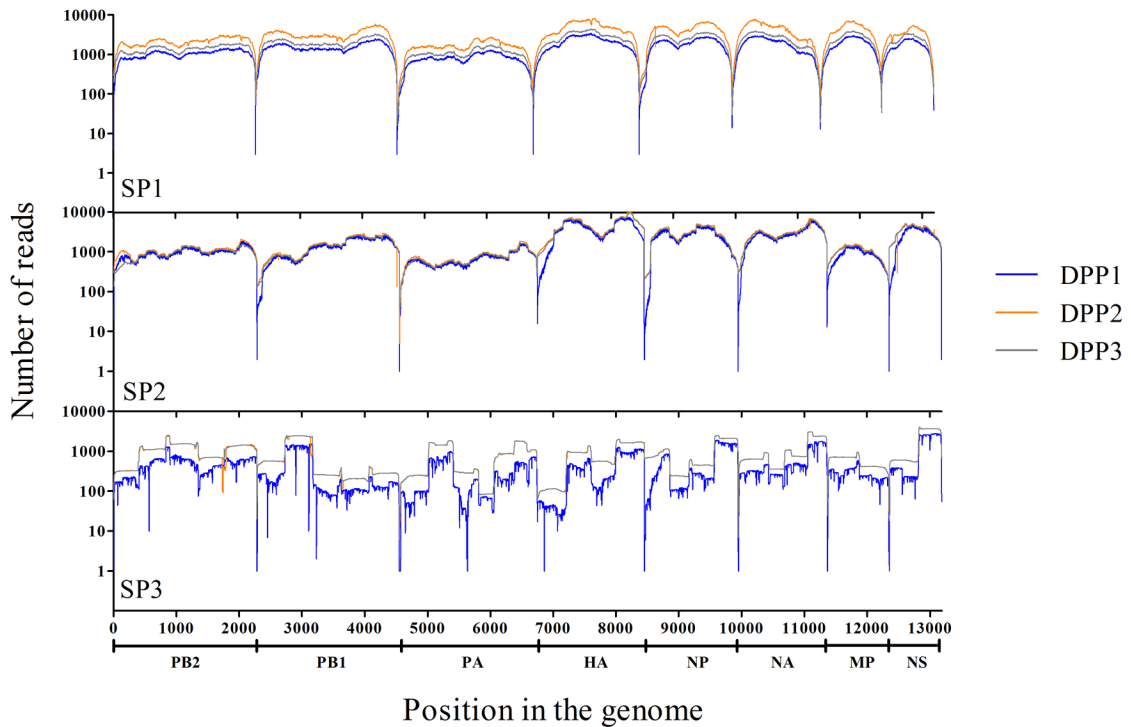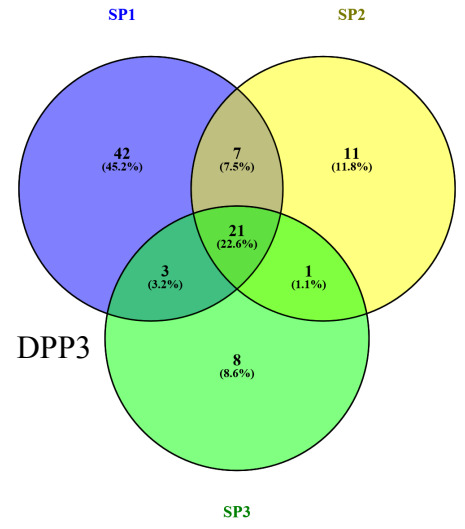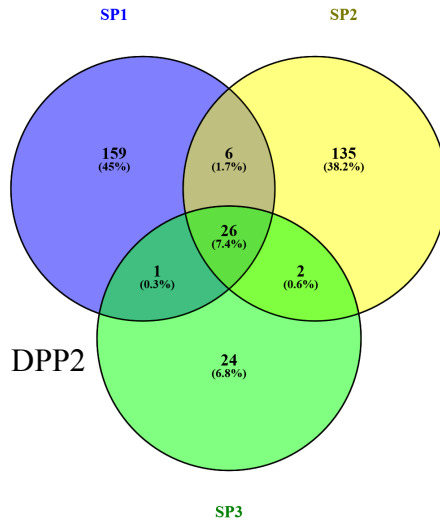
795    (purple) and >10% (green)

796

797

798

**1. Sharing virus stocks**

UKDD

DETU

NLCH

**2. Sample preparation, sequencing and raw data trimming**

Institution 1  -  SP1 (Illumina MiSeq)  -  DPP1 (FluSeqID)

Institution 2  -  SP2 (Illumina MiSeq)  -  DPP2 (Genome Sequencer & Geneious software)

Institution 3  -  SP3 (Roche 454)       -  DPP3 (CLC Bio)

**3. Sequence data sharing**

Data Hub

**4. Sequence data alignment and generate consensus sequences**

SP1, SP2, SP3 data in DPP1

SP1, SP2, SP3 data in DPP2

SP1, SP2, SP3 data in DPP3

**5. One consensus**

Per virus

Per gene segment

**6a. Sequence data alignment and mSNV analysis**

SP1, SP2, SP3 data in DPP1

SP1, SP2, SP3 data in DPP2

SP1, SP2, SP3 data in DPP3

**6b. Sequence data alignment**

SP1, SP2, SP3 data in DPP1

SP1, SP2, SP3 data in DPP2

SP1, SP2, SP3 data in DPP3

**6c. mSNV analysis**

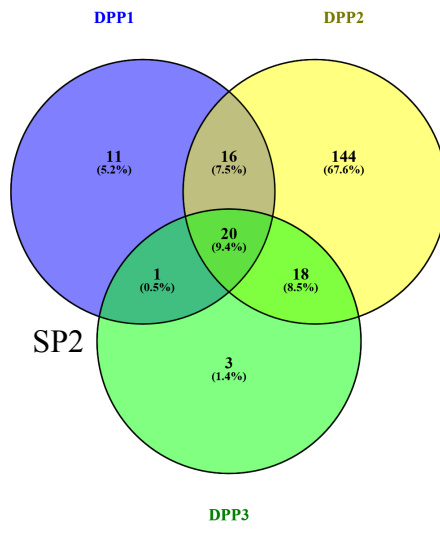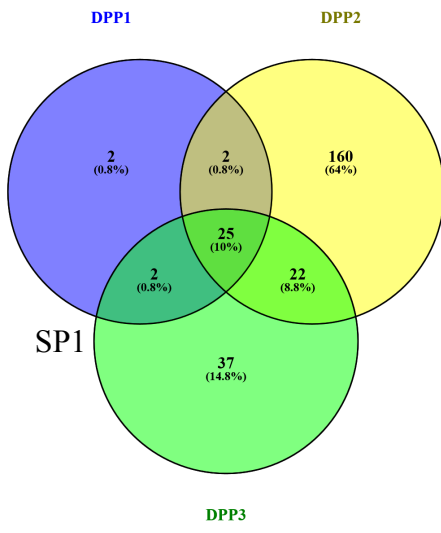DPP3

**7. Data comparison**

**A: SP derived differences: SP datasets anaysed per DPP**

**B: DPP derived differences: DPP anayses results per SP**

| Gene segment | Set | Sense | Primer Sequence |
|---|---|---|---|
| PB2 | 1 | 3-Forward | CGAAAGCAGGTCAAATATATTC |
| | | 521-Reverse | TCCATGATGACATCTTGTGCTTC |
| | 2 | 428-Forward | CATGGAACCTTCGGTCCCGTTCA |
| | | 931-Reverse | ATCCACAGCTTGTTCCTCAGTTGG |
| | 3 | 855-Forward | AGCAACGGTATCAGCGGATCCA |
| | | 1403-Reverse | CCATGACATTATCAATGGGTTC |
| | 4 | 1315-Forward | CCCATGCATCAACTCCTGAGACA |
| | | 1820-Reverse | GTTCTCACAAATCCACTGTATTG |
| | 5 | 1759-Forward | GAACCGTTCCAATCCTTGGTACCT |
| | | 2341-Reverse | AGTAGAAACAAGGTCGTTT |
| PB1 | 1 | 3-Forward | RAAAGCAGGCAAACCAYTTGAATG |
| | | 538-Reverse | CCATCACATCCTTGAGGAAATC |
| | 2 | 445-Forward | ACYGCTTTGGCCAACACTATAGA |
| | | 944-Reverse | GTATTGTCYCCATGAATTGTAAAGG |
| | 3 | 877-Forward | GTCCTCAGGAACATGATGACTAACTCAC |
| | | 1403-Reverse | ATTCCCTCATGATTCGGTGC |
| | 4 | 1319-Forward | CCAAAACCACATATTGGTGGGACGG |
| | | 1892-Reverse | CTGCCCTGGTARTCTTCATCCATC |
| | 5 | 1782-Forward | GGCAGGACTGTTGGTTTCAGATGG |
| | | 2326-Reverse | TTTTTTCAYGAAGGACAAGC |
| PA | 1 | 3-Forward | CRAAAGCAGGTACTGATYC |
| | | 607-Reverse | CGGATTGACGAAAGGAATCCCA |
| | 2 | 452-Forward | CACACATTCACATATTCTCATTCAC |
| | | 897-Reverse | GCTTAATTTAAGYGCATCCATTCAC |
| | 3 | 731-Forward | GAGGGCAAGCTTTCTCAAATGTC |
| | | 1305-Reverse | TTCATCAAGTTCAATCCAACTGA |
| | 4 | 1168-Forward | GAGGACTGCAAAGATGTTAGCGA |
| | | 1646-Reverse | CAGTACTTTTCCCACTTGTGTGG |
| | 5 | 1490-Forward | GCAGAACCAAAGAAGGAAGACGG |
| | | 2072-Reverse | GATCGAAGGTCCCAGGTTCCAGG |
| | 6 | 1816-Forward | GCCGAGTCTTCTGTCAAAGAGAA |
| | | 2233-Reverse | AGTAGAAACAAGGTACYTTTT |
| HA | 1 | 5-Forward | AAAGCAGGGGTHYDATCTGTC |
| | | 570-Reverse | TTGTARCTYCTCTTTATBGTBGG |
| | 2 | 465-Forward | GRGTRAGCKCAGCATGTCC |
| | | 917-Reverse | GDGTTTGRCACTTGGTGTTGC |
| | 3 | 803-Forward | AGTAATGGRAATTTCATTGCYCC |
| | | 1378-Reverse | ATTYTCCATKAGAACYAGRAGTTC |
| | 4 | 1247-Forward | ACTCARTTTGARGCHGTTGG |

| | | | |
|---|---|---|---|
| | | 1789-Reverse | AGTAGAAACAAGGGTGTTTT |
| NP | 1 | 1-Forward | AGCRAAAGCAGGGTDKATA |
| | | 482-Reverse | GCATCATTYAGRTTKGAATGCC |
| | 2 | 239-Forward | GAATGGTNCTCTCTGCVTTTG |
| | | 838-Reverse | TGAGTGCAGACCGHGCCAG |
| | 3 | 729-Forward | RAAATTYCAAACAGCAGCAC |
| | | 1266-Reverse | CTKATYTGYCCTGCVGATGC |
| | 4 | 1132-Forward | GTTCAAATTGCTTCAAATG |
| | | 1565-Reverse | AGTAGAAACAAGGGTATTTT |
| NA | 1 | 3-Forward | CRAAAGCAGGAGTTYAAAATG |
| | | 531-Reverse | GGCTTGATATACATTGGGTGATTG |
| | 2 | 400-Forward | TGCAGGACTTTCTTCCTCACTCA |
| | | 900-Reverse | GTTGTCTCTACACACGCATTCCAC |
| | 3 | 731-Forward | ATTGGGTAATGACTGACGGTCC |
| | | 1237-Reverse | AAGACCCACTGTATCCCGACCA |
| | 4 | 1103-Forward | GGACAATTAGTCGAACCTCCAGA |
| | | 1460-Reverse | AGTAGAAACAAGGAGTTTTT |
| MA | 1 | 5-Forward | AAAGCAGKTAGATRTTGAAARATG |
| | | 564-Reverse | ACCATTCTGTTYTCATGYCTG |
| | 2 | 461-Forward | TAKTRTGTGCCACTTGTGAGC |
| | | 1023-Reverse | AGTAGAAACAAGGTARKTTTT |
| NS | 1 | 3-Forward | CRAAAGCAGGGTGACAAAVAC |
| | | 547-Reverse | CCAATTGCAWTYTTGACATCCTC |
| | 2 | 453-Forward | AGAGCTTTCACRGAAGAAGGAGCA |
| | | 888-Reverse | AGTAGAAMCAAGGGTGTTTT |

| Virus | Segment | Covering positions from start codon ATG | Consensus sequence |
|---|---|---|---|
| NLCH | PB2 | 1-2280 | **ATG**GAGAGAATAAAAGAACTAAGAGATCTAATGTCTCAATCCCGC ACTCGCGAGATACTAACAAAAACCACTGTGGACCATATGGCCATA ATCAAGAAATACACATCAGGAAGACAAGAGAAGAACCCTGCTCTC AGAATGAAATGGATGATGGCAATGAAATATCCAATCACAGCAGAC AAGAGAATAATGGAAATGATTCCTGAAAGAAATGAACAAGGCCA GACGCTTTGGAGCAAGACAAATGATGCTGGATCAGACAGAGTGAT GGTGTCTCCCCTAGCTGTAACTTGGTGGAATAGAAATGGACCGAC AGCAAGTACAGTCCATTATCCAAAGGTCTACAAAACATACTTTGA GAAGGTTGAAAGGTTAAAGCATGGAACCTTCGGTCCCGTTCACTTC CGAAACCAAATTAAAATACGCCGCCGAGTTGACATAAACCCAGGC CACGCAGATCTCAGTGCCAAAGAAGCACAAGATGTCATCATGGAG GTCGTTTTCCCAAATGAAGTGGGAGCTAGAATATTGACATCAGAG TCACAATTGACAATAACGAAAGAGAAAAAGAAGAACTCCAGGA TTGCAAGATTGCTCCTTTAATGGTGGCATACATGTTGGAAAGAGAA CTGGTCCGCAAAACCAGATTCCTACCAGTAGCAGGTGGGACAAGC AGTGTGTACATTGAGGTACTGCACCTGACCCAAGGGACCTGCTGG GAACAGATGTACACTCCAGGCGGAGAAGTGAGAAATGACGATGTT GACCAGAGTTTGATCATCGCGGCCAGAAACATTGTTAGGAGAGCA ACGGTATCAGCGGATCCACTGGCATCATTATTGGAGATGTGCCAC AGCACACAAATTGGTGGGACAAGGATGGTGGATATCCTTAGGCAA AATCCAACTGAGGAACAAGCTGTGGATATATGCAAAGCAGCAATG GGTTTAAGGATTAGTTCATCCTTTAGCTTTGGAGGATTCACCTTCA AAAGAACTAGTGGTTCATCCATTAGAAAGGAAGAGGAAGTGCTTA CAGGCAACCTCCAAACATTGAAAATAAGAGTACATGAGGGGTAT GAGGAGTTCACAATGGTTGGGCGAAGAGCAACAGCCATTCTAAGG AAAGCAACTAGAAGGCTGATTCAGTTGATAGTAAGTGGAAGAGAC GAACAATCAATCGCTGAAGCAATCATCGTAGCCATGGTGTTCTCAC AGGAGGATTGCATGATAAAGGCAGTCCGAGGCGATCTAAATTTT GTGAACAGAGCAAACCAAAGATTGAACCCCATGCATCAACTCCTG AGACACTTCCAAAAAGATGCAAAAGTGCTGTTTCAAAAATTGGGGG ATCGAACCCATTGATAATGTCATGGGGATGATTGGAATATTGCCTG ACATGACTCCAAGCACAGAGATGTCACTAAGAGGAGTAAGAGTT AGTAAAATGGGAGTAGATGAATATTCCAGCACTGAGAGAGTGGTT GTAAGCATTGACCGTTTCTTGCGGGTTCGAGATCAGCAGGGGAAC GTACTCCTATCTCCCGAAGAAGTCAGCGAAACACTGGGAACAGAA AAATTAACAATAACATATTCATCATCAATGATGTGGGAAATCAAT GGTCCTGAGTCAGTGCTGGTCAACACCTATCAATGGATCATCAGA AATTGGGAGATTGTGAAGATTCAATGGTCTCAAGACCCCACGATG CTGTACAATAAGGTGGAGTTTGAACCGTTCCAATCCTTGGTACCTA AAGCTGCCAGAGGCCAATACAGTGGATTTGTGAGAACACTGTTC CAACAAATGCGTGACGTATTGGGGACATTTGATACTATTCAGATA ATAAAGCTGTTACCGTTTGCAGCAGCCCCACCGGAGCATAGCAGA ATGCAATTTTCTTCCCTGACCGTGAATGTAAGGGGCTCGGGAATGA GAATACTCGTAAGGGGTAACTCCCCTGTGTTCAACTACAATAAG GCAACCAAAAGGCTTGCCGTCCTTGGAAAGGACGCAGGTGCATTA ACAGAGGATCCAGATGAGGGGACAACAGGAGTGGAATCTGCAGT GCTGAGGGGGTTCCTAATTCTGGGCAGGGAGGACAGAAGATATGG ACCAGCACTAAGCATCAATGAACTGAGCAATCTTGCGAAAGGGGA GAAAGCCAATGTGCTGATAGGGCAAGGAGACGTGGTGCTGGTAAT GAAACGGAAACGGGACTCTAGCATACTTACTGACAGCCAGACAGC GACCAAAAGAATTCGGATGGTCATCAATTAG |
| NLCH | PB1 | 1-2277 | **ATG**GATGTCAACCCGACTCTACTCTTCTTGAAAGTGCCAGCGCAA |

| | | | |
|---|---|---|---|
| | | | AATGCTATAAGTACCACATTCCCCTATACTGGAGATCCTCCATACA<br>GCCATGGAACAGGAACAGGATACACCATGGACACAGTCAACAGA<br>ACGCATCAATACTCAGAAAAGGGAAAGTGGACAAAAAACACCGA<br>GACTGGAGCACCCCAACTCAACCCAATTGATGGACCATTACCTGA<br>GGATAACGAGCCAAGCGGATATGCACAAACGGATTGTGTGTTGGA<br>AGCAATGGCTTTCCTTGAAGAGTCCCACCCAGGGATCTTTGAAAAC<br>TCATGTCTTGAAACAATGGAAATTGTTCAACAAACAAGAGTGGAC<br>AAACTGACCCAAGGTCGTCAGACCTATGACTGGACATTGAATAGA<br>AACCAGCCGGCTGCAACTGCTTTAGCCAACACTATAGAAGTCTTCA<br>GATCGAACGGTCTAACAGCCAATGAGTCAGGGAGACTGATAGATT<br>TCCTCAAAGATGTGATGGAGTCAATGGACAAAGAAGAAATGGAA<br>ATAACAACACATTTCCAAAGAAAGAGAAGAGTAAGAGACAATAT<br>GACCAAGAAAATGGTCACACAAAGAACAATAGGGAAGAAAAAAC<br>AGAGACTGAACAAGAAGAACTACTTGGTAAGGGCACTGACACTGA<br>ACACAATGACAAAAGATGCAGAAAGAGGCAAGTTGAAGAGGCGG<br>GCAATTGCAACACCCGGGATGCAAATCAGAGGGTTCGTGTACTTT<br>GTCGAAACATTAGCGAGGAGCATCTGCGAGAAACTTGAGCAATCT<br>GGGCTCCCTGTTGGAGGAAATGAAAAAAAGGCTAAGTTGGCAAAT<br>GTCGTGAGAAAGATGATGACTAACTCACAAGACACAGAGCTATCC<br>TTTACAATTACTGGAGACAATACCAAGTGGAACGAGAATCAGAAT<br>CCTCGGATTTTTTTGGCAATGATAACATATATCACAAGAAATCAAC<br>CTGAGTGGTTTAGAAATGTGTTAAGTATTGCCCCTATAATGTTCTC<br>AAACAAAATGGCAAGGTTAGGGAAAGGATACATGTTCGAAAGTA<br>AGAGCATGAAGCTACGGACACAAATACCAGCAGAAATGCTTGCAA<br>CCATTGACCTGAAATATTTCAACGAATCGACAAGAAAGAAAATTG<br>AGAAAATAAGGCCTCTCCTAATAGAAGGGACAGCCTCGTTGAGTC<br>CTGGAATGATGATGGGCATGTTCAACATGCTGAGTACAGTCTTGG<br>GAGTATCAATTCTAAATCTTGGCCAAAAGAGGTACACCAAAACCA<br>CATACTGGTGGGACGGACTCCAATCCTCTGATGATTTCGCTCTCAT<br>AGTAAATGCACCGAATCATGAGGGAATACAGGCAGGAGTGGACA<br>GGTTCTATAGGACTTGTAAATTGGTTGGGATCAATATGAGTAAAA<br>AGAAATCCTATATAAATCGGACAGGAACATTTGAATTCACAAGCT<br>TTTTCTACCGTTATGGGTTTGTAGCCAACTTCAGCATGGAGCTGCC<br>CAGCTTTGGAGTTTCTGGGATCAATGAATCGGCTGACATGAGCATT<br>GGAGTTACAGTAATAAAGAATAACATGATAAACAACGATCTTGGA<br>CCAGCAACAGCTCAAATGGCTCTTCAGCTATTTATCAAGGACTACA<br>GATATACATATCGATGCCACAGGGGTGATACACAAATACAAACAA<br>GGAGATCATTCGAGCTAAAGAAGCTGTGGGAGCAGACCCGTTCAA<br>AGGCAGGACTGTTGGTTTCAGATGGAGGCCCAAACTTATACAATA<br>TACGGAATCTCCACATCCCAGAGGTCTGCTTGAAGTGGGAACTGA<br>TGGATGAAGATTACCAGGGTAGACTTTGTAATCCCCTGAACCCCTT<br>TGTCAGTCATAAGGAAATTGAATCCGTAAACAATGCTGTAGTGAT<br>GCCAGCCCATGGTCCGGCCAAAAGCATGGAATATGATGCTGTTGC<br>GACCACACACTCATGGGTCCCTAAGAGGAACCGTTCCATTCTGAAT<br>ACCAGTCAAAGAGGAATCCTTGAGGATGAACAGATGTATCAGAAG<br>TGCTGCAATCTATTTGAAAAATTCTTCCCTAGTAGCTCATACAGGA<br>GGCCAGTTGGAATCTCCAGTATGGTGGAGGCCATGGTGTCTAGGG<br>CCCGAATTGATGCACGGATTGACTTCGAGTCTGGTAGGATTAAGA<br>AGGAAGAGTTTGCTGAGATCATGAAGATCTGTTCCACCATTGAAG<br>AGATCAGACGGCAAAAACAGTGA |
| NLCH | PA | -6-2190 | TCCAAA**ATG**GAAGACTTTGTGCGACAATGCTTCAATCCAATGATC<br>ATCGAGCTTGCGGAAAAGACAATGAAAGAATATGGGGAAAATCC<br>AAAAATCGAAACGAACAAATTCGCTGCAATATGCACTCACTTAGA<br>GGTCTGTTTCATGTATTCGGATTTCCACTTTATTGATGAACGTGGT<br>AAATCAATAATTGTAGAATCTGGCGATCCGAATGCATTATTGAAA<br>CACCGATTTGAGATAATTGAAGGGAGGGACCGAACGATGGCTTGG |

| | | | |
|---|---|---|---|
| | | | ACAGTGGTAAATAGTATCTGCAACACCACAGGAGTCGATAAGCCT<br>AAATTCCTCCCAGATTTGTATGATTACAAGGAGAACCGATTCATT<br>GAAATTGGAGTGACAAGGAGGGAAGTTCACACATACTACCTAGAA<br>AAGGCAAATAAGATAAAATCAGAGAAGACACACATTCACATATTC<br>TCATTCACTGGGGAGGAGATGGCCACCAAAGCTGATTATATCCTTG<br>ATGAAGAGAGCAGGGCAAGGATCAAAACCAGGTTGTTCACTATC<br>AGGCAAGAAATGGCCAATAGGGGTCTGTGGGATTCCTTTCGTCAA<br>TCTGAGAGAGGCGAAGAGACAATTGAAGAAAGGTTTGAAATCACA<br>GGAACCATGCGCAGGCTTGCCGACCAAAGTCTCCCACCGAATTTCT<br>CCAGCCTTGAAAATTTTAGAGCCTATGTGGATGGATTCAAACCG<br>AACGGCTGCCTTGAGGGCAAGCTTTCTCAAATGTCAAAAGAAGTG<br>AACGCCAGAATTGAGCCATTCATGAAGAAAACACCACGCCCTCTC<br>AGATTACCTGATGGTCCTCCTTGCTCTCAGCGGTCGAAATTCTTAC<br>TGATGGATGCTCTTAAATTGAGCATCGAAGACCCAAGCCATGAG<br>GGAGAAGGTATACCGCTATATGATGCAATCAAATGCATGAAGACG<br>TTTTTTGGTTGGAAAGAGCCCAACATTGTAAAACCACATGTAAAA<br>GGCATAAATCCCAACTATCTCTTGGCTTGGAAGCAGGTGCTGGTAG<br>AACTCCAAGACATTGAAAATGAAGAGAAAATCCCAAAAACAAAA<br>AACATGAAGAAAACAAGCCAACTAAAATGGGCACTCGGTGAGAA<br>TATGGCACCTGAAAAAGTGGACTTTGAGGACTGCAAAGATGTTAG<br>CGATCTAAGACAGTATGACAGTGATGAACCAGAGCCCAGATCATT<br>ATCAAGCTGGATCCAGAGCGAATTCAACAAAGCATGCGAATTGAC<br>AGATTCGAGTTGGATTGAACTTGATGAAATAGGAGAAGATGTTGC<br>TCCAATTGAGCACATTGCGAGTATGAGAAGAAACTACTTCACAGC<br>GGAAGTGTCTCATTGCAGGGCTACTGAATATATAATGAAAGGAGT<br>TTATATAAATACAGCCCTGTTGAATTCATCCTGTGCAGCCATGGAT<br>GACTTCCAATTGATTCCAATGATAAGCAAGTGCAGAACCAAAGAA<br>GGAAGACGGAAGACAAATCTATATGGGTTCATTATAAAAGGAAGA<br>TCCCATTTGAGGAATGATACCGATGTGGTAAATTTTGTGAGCATGG<br>AGTTCTCTCTTACTGACCCGAGGCTGGAACCACACAAGTGGGAA<br>AAGTACTGTGTTCTCGAAATAGGAGACATGCTCCTACGAACTGCA<br>ATAGGCCAAGTATCAAGACCCATGTTTCTTTATGTAAGGACCAATG<br>GGACTTCCAAGATCAAGATGAAATGGGGCATGGAGATGAGGCGAT<br>GCCTTCTTCAATCCCTCCAACAAATTGAGAGCATGATTGAGGCA<br>GAGTCTTCTGTCAAAGAGAAGGACATGACCAAGGAATTCTTTGAA<br>AATAAATCAGAAACGTGGCCAATTGGGGAATCACCTAAGGGGGTG<br>GAGGAAAGCTCTATTGGGAAAGTGTGTAGAACATTACTAGCAAAA<br>TCTGTATTCAACAGCCTATATGCATCTCCACAACTTGAGGGGTTT<br>TCAGCTGAGTCGAGAAAGTTACTTCTCATTGTTCAGGCATTTAGGG<br>ACAACCTGGAACCTGGGACCTTCGATCTTGGGGGGCTATATGAAG<br>CAATTGAGGAGTGCCTGATTAATGATCCCTGGGTTTTGCTTAATGC<br>ATCTTGGTTCAACTCCTTCCTTACACATGCACTGAAATAGTTG<br>TGGCAATGCTACTATTTGCTATCCATACTGTCCAAA |
| NLCH | HA | 7-1704 | AAAATAGTGCTTCTTCTTGCAGTGGTTAGCCTTGTTAAAAGTGATC<br>AGATTTGCATTGGTTACCATGCAAACAACTCAACAAAACAGGTTG<br>ACACAATAATGGAAAAAAACGTCACTGTTACACATGCCCAAGACA<br>TACTGGAAAAGACACACAACGGGAAGCTCTGCGATCTTAATGGA<br>GTGAAGCCCCTGATTCTAAAGGATTGTAGCGTAGCTGGGTGGCTCC<br>TTGGAAATCCAATGTGCGACGAGTTCATCAGGGTGCCGGAATGGT<br>CTTACATCGTGGAGAGGGCTAACCCAGCCAACGACCTCTGTTACCC<br>AGGGACCCTCAATGACTATGAGGAACTGAAACACCTACTGAGC<br>AGAATAAATCATTTTGAGAAAACTCTGATCATCCCCAAGAGTTCTT<br>GGCCCAATCATGAAACATCATTAGGGGTGAGCGCAGCATGTCCAT<br>ACCAGGGAGCATCCTCATTTTTCAGAAATGTGGTATGGCTCATCAA<br>AAAGAACGATGCATACCCGACAATAAAGATAAGCTACAATAAT<br>ACCAATCGGGAAGATCTTTTGATACTGTGGGGGATTCATCATCCCA |

| | | | |
|---|---|---|---|
| | | | ACAATGCAGAAGAGCAGACAAATCTCTATAAAAACCCAGACACTT ATGTTTCCGTTGGGACATCAACATTAAACCAGAGATTGGTGCCAA AAATAGCTACTAGATCCCAAGTAAACGGGCAACGTGGAAGAATG GATTTCTTCTGGACAATTTTAAAACCGAATGATGCAATCCACTTTG AGAGTAATGGAAATTTCATTGCTCCAGAATATGCCTACAAAATTGT CAAGAAAGGGGACTCAACAATTATGAAAAGTGAAGTGGAGTATG GCCACTGCAACACCAAATGTCAAACCCCAATAGGGGCGATAAAC TCTAGCATGCCATTCCACAATATACACCCTCTCACCATCGGGGAAT GCCCCAAATACGTGAAGTCAAACAAATTAGTCCTTGCGACTGGGC TCAGAAATAGTCCTCTAAGGGAAAGAAGAAGAAAAAGAGGACTA TTTGGAGCTATAGCAGGGTTTATAGAGGGAGGATGGCAGGGAATG GTAGACGGTTGGTATGGGTACCACCATAGCAATGAGCAGGGGAGT GGGTACGCTGCAGACAAAGAATCCACCCAAAAGGCAGTAGATGG AGTTACCAATAAGGTCAACTCAATCATTGACAAAATGAACACTCA ATTTGAGGCCGTTGGAAGGGAATTTAATAATTTAGAAAGGAGAAT AGAGAATCTAAACAAGAAAATGGAAGACGGATTCCTAGATGTCTG GACTTATAATGCTGAACTTTTAGTTCTCATGGAAAATGAGAGAACT CTAGATTTCCATGACTCAAATGTCAAGAACCTTTACGACAAAGTCC GACTACAGCTTAGGGATAATGCAAAGGAGCTGGGTAATGGTTGT TTCGAGTTCTATCACAAATGTGATAACGAATGTATGGGAAGCGTA AGAAATGGGACGTATGACTACCCTAAGTATTCAGAAGAAGCAAGG TTAAAAAGAGAAGAAATAAGCGGAGTGAAATTAGAATCAATAGG AACTTACCAAATACTGTCAATTTATTCAACAGTGGCGAGTTCCCTA GCACTGGCAATCATAGTGGCTGGTCTATCTTTATGGATGTGCTCTA ATGGGTCGCTACAATGCAGAATTTGCATCTAA |
| NLCH | NP | 1-1497 | **ATG**GCGTCTCAAGGCACCAAACGATCTTATGAACAGATGGAAACT GGTGGAGAACGCCAGAATGCCACTGAAATCAGAGCATCTGTTGGA AGAATGGTTGGTGGAATTGGAAGGTTTTATATACAGATGTGCACT GAACTCAAACTCAGCAATTATGAGGGGAGACTGATCCAGAACAGC ATAACAATAGAAAGAATGGTTCTCTCTGCATTTGATGAAAGGAGG AACAAGTACCTGGAAGAACATCCCAGTGCGGGGAAGGACCCAAA GAAAACTGGAGGTCCAATCTACAGAAGAAGAGACGGAAAGTGGA TGAGGGAGCTGATTCTGTATGACAAAGAAGAGATCAGAAGGATCT GGCGTCAAGCAAATAATGGAGAAGATGCAACTGCTGGTCTCACCC ATCTGATGATCTGGCACTCCAACCTGAATGATGCCACATATCAGAG GACAAGGGCTCTCGTGCGCACTGGAATGGATCCCAGAATGTGCTC TCTGATGCAAGGATCAACTCTCCCAAGAAGGTCTGGAGCTGCTGG TGCAGCAGTAAAAGGGGTCGGAACAATGGTAATGGAATTGATTCG AATGATAAAGCGAGGGATTAATGATCGGAATTTCTGGAGAGGCGA AAATGGAAGAAGGACAAGGATTGCCTATGAGAGAATGTGCAACAT CCTCAAAGGGAAATTTCAAACAGCAGCACAAAGAGCAATGATGGA TCAAGTGCGAGAAAGCAGGAATCCTGGGAATGCTGAAATTGAAGA TCTCATTTTTCTGGCACGGTCTGCACTCATCCTGAGAGGATCAGTG GCCCACAAGTCTTGTCTGCCTGCTTGTGTTTACGGACTTGCTGTGG CCAGTGGATATGACTTTGAGAGAGAAGGATACTCTCTGGTTGGA ATAGACCCTTTCCGTCTGCTTCAAAACAGCCAGGTCTTCAGTCTCA TTAGACCAAATGAAAACCCAGCACATAAAAGTCAGTTGGTATGGA TGGCATGCCATTCAGCAGCGTTTGAGGACCTGAGAGTATCAAGTTT CATCAGAGGGACAAGAGTGGTCCCAAGAGGACAACTATCCACCAG AGGAGTTCAAATTGCATCAAATGAAAACATGGAAACAATGGACTC CAGCACTCTTGAATTGAGAAGCAGATACTGGGCTATAAGAACCAG GAGTGGAGGAAACACCAACCAACAGAGAGCTTCTGCAGGACAAA TCAGCGTACAACCCACCTTCTCAGTACAGAGAAATCTTCCCTTT GAAAGAACGACCATCATGGCGGCATTTACAGGGAACACTGAAGGC AGGACCTCTGACATGAGGACTGAGATCATAAGAATGATGGAAAGT GCCAAACCAGAAGATGTGTCCTTCCAGGGGCGGGGAGTCTTCGAG |

| | | | |
|---|---|---|---|
| | | | CTCTCGGACGAAAAGGCAACGAACCCGATCGTGCCTTCCTTTGAC<br>ATGAGCAACGAAGGATCTTATTTCTTCGGAGACAGTGCAGAGGAG<br>TATGACAATTAA |
| NLCH | NA | 4-1419 | AATCCAAATCAGAAAATAGTAACCATTGGCTCCATTTCATTAGGGT<br>TGGTTGTATTCAATGTTCTACTGCATGCTGTGAGCATCATATTAAC<br>AGTGTTAGCCCTGGGGAAGAGTGAAAACAATGGAATCTGCAATGG<br>AACTGTAGTGAGAGAATACAATGAAACAGTTAGAATAGAGAAA<br>GTGACTCAATGGTACAATACTAGCGTAGTCGAATATGTACCGCATT<br>GGAATGAGGGCACTTATATAAATAACACCGAACCAATATGTGATG<br>TCAAGGGCTTTGCACCTTTTTCCAAGGACAACGGGATAAGAGTTG<br>GCTCCAGGGGACATATTTTTGTCATAAGAGAGCCTTTCGTCTCT<br>TGTTCACCTGTAGAGTGCAGGACTTTCTTCCTCACTCAGGGATCTC<br>TACTCAATGACAAACACTCAAATGGAACAATGAAGGATAGAAGCC<br>CATTCAGAACTCTCATGAGTGTCGAAGTGGGCCAATCACCCAATGT<br>ATATCAAGCCAGGTTTGAAGCTGTGGCATGGTCAGCAACAGCC<br>TGTCATGATGGTAAGAAGTGGATGACGATTGGTGTAACAGGGCCA<br>GATTCTAAAGCAGTAGCAGTAGTTCATTACGGAGGGGTGCCTACT<br>GACGTTGTTAACTCCTGGGCAGGAGATATATTAAGAACTCAGGAG<br>TCATCTTGTACTTGCATTCAAGGTAATTGTTATTGGGTAATGACT<br>GACGGTCCTGCCAATAGACAGGCGCAGTATAGAATATACAAAGCA<br>AATCAAGGCAAAATAATTGGCCGAACAGATGTTAGCTTTAGTGGA<br>GGACATATTGAGGAATGTTCTTGTTATCCAAATGATGGTAAAGTGG<br>AATGCGTGTGTAGAGACAACTGGACGGGAACTAACAGGCCTGTA<br>CTAATTATTTCGCCTGATCTCTCTTACAGGGTTGGGTATTTATGTGC<br>AGGGTTGCCCAGTGACACTCCAAGAGGGGAAGATACTCAATTTGT<br>CGGTTCATGCACTAGTCCCATGGGAAATCAGGGATATGGCGTAAA<br>AGGGTTCGGGTTTCGACAGGGAACTGATGTGTGGGTGGGGCGG<br>ACAATTAGTCGAACCTCCAGATCAGGATTTGAAATAATAAGGATA<br>AAGAATGGTTGGACGCAAACAAGCAAAGAACAGATTAGAAGACA<br>GGTGGTTGTTGATAACTCGAATTGGTCGGGATACAGTGGGTCTTTC<br>ACTTTACCAGTAGAATTGTCTGGGAGGGAATGTTTGGTTCCCTGT<br>TTTTGGGTCGAAATGATCAGAGGTAGGCCAGAAGAGAGAACAATC<br>TGGACCTCTAGTAGCTCCATTGTAATGTGTGGAGTTGATTATGAAA<br>TTGCCGATTGGTCATGGCACGATGGAGCTATTCTTCCCTTTGACAT<br>CGATAAGACGTAATTTACG |
| NLCH | MP | -5-982 | GAAAG**ATG**AGTCTTCTAACCGAGGTCGAAACGTACGTTCTCTCTAT<br>CATCCCGTCAGGCCCCCTCAAAGCCGAGATCGCGCAGAGACTTGA<br>AGATGTCTTTGCAGGGAAAAACACCGATCTCGAGGCTCTCATGGA<br>GTGGCTAAAGACAAGACCAATCCTGTCACCTCTGACTAAAGGGA<br>TTTTGGGATTTGTGTTCACGCTCACCGTGCCCAGTGAGCGAGGACT<br>GCAGCGTAGACGCTTCGTCCAGAATGCCCTAAATGGAAACGGGGA<br>TCCAAATAATATGGATAAGGCAGTTAAGCTATATAAGAAGCTGAA<br>AAGAGAGATAACATTCCATGGGGCTAAGGAGGTCGCACTTAGCT<br>ACTCAACCGGTGCACTTGCCAGCTGCATGGGTCTCATATACAACAG<br>GATGGGAACGGTGACTACAGAAGTGGCTTTTGGCCTAGTGTGTGC<br>CACTTGTGAGCAGATTGCAGATTCACAGCATCGGTCCCACAGACA<br>GATGGCAACCATCACCAACCCATTAATCAGACATGAGAACAGAA<br>TGGTGCTGGCCAGCACTACAGCTAAGGCCATGGAGCAGATGGCAG<br>GATCAAGCGAGCAGGCATCAGAAGCCATGGAGGTTGCTAATCAGG<br>CCAGGCAGATGGTACAGGCAATGAGGACAATTGGGACTCATCCTA<br>ACTCTAGTGCTGGTCTGAGAGATAATCTTCTTGAAAATTTGCAGG<br>CCTACCAGAACCGAATGGGAGTGCAGATGCAGCGATTCAAGTGAT<br>CCTCTTGTTGTTGCCGCAAATATCATTGGGATCCTGCACTTGATATT<br>GTGGATCCTTGATCGTCTTTTCTTCAAATGCATTTATCGTCGCCTTA<br>AATACGGTTTGAAAATAGGGCCTTCTACGGAAGGGGTACCT<br>GAGTCTATGAGGGAAGAGTACCGGCAGGAACAGCAGAGTGCTGTG |

| | | | GATGTTGACGATGGTCATTTTGTCAACATAGAATTGGAGTAA |
|---|---|---|---|
| NLCH | NS | 1-838 | **ATG**GACTCCAACACTGTGTCAAGCTTTCAGGTAGACTGCTTTCTTT<br>GGCATGTCCGCAAACGATTTGCAGACCAAGAACTGGGTGATGCCC<br>CATTCCTTGACCGGCTTCGCCGAGACCAGAAGTCCCTAAGAGGAA<br>GAGGCAGCACTCTTGGTCTGGACATCGAGACAGCTACTCGTGCG<br>GGAAAGCAAATATTGGAGCGGATTCTGGGGGAAGAATCTGATGAA<br>GCACTTAAAATGAATATTGCTTCTGTACCGACTTCACGCTACCTAA<br>CTGACATGACTCTTGAAGAAATGTCAAGAGACTGGTTCATGCTCAT<br>GCCCAAGCAGAAAGTAGCAGGTTCTCTCTGCATCAAAATGGACCA<br>GGCAATAATGGATAAAACCATCATACTGAAAGCAAACTTCAGTGT<br>GATTTTTGATCGGCTGGAAACCCTAATATTACTTAGAGCTTTCACA<br>GAAGAAGGAGCAATTGTGGGAGAAATCTCACCATTACCTTCTCTTC<br>CAGGACATACTGATGAGGATGTCAAAATTGCAATTGGGGTCCTCA<br>TCGGAGGGCTTGAATGGAATGATAACACAGTTCGAGTCTCTGAAA<br>CTCTACAGAGATTCACTTGGAGAAGCAGTAATGAGGATGGGAGAC<br>CTTCACTCCCTTCAAAACAGAAACGGAAAATGGCGAGAACAATTG<br>AGTCAGAAGTTCGAGGAAATAAGGATGGCTGATTGAGGAAATGCGA<br>CATAGATTGAAGATCACAGAGAACAGCTTCGAACAAATAACGTTT<br>ATGCAAGCTTTACAACTATTGCTTGAAGTGGAGCAAGAGATAAGA<br>ACCTTCTCGTTTCAGCTTATTTAA |
| DETU | PB2 | 1-2287 | **ATG**GAGAGAATAAAAGAACTAAGAGATCTAATGTCTCAATCCCGC<br>ACTCGCGAGATACTAACAAAAACCACTGTGGACCATATGGCCATA<br>ATCAAGAAATACACATCAGGAAGACAAGAGAAGAACCCTGCTCTC<br>AGAATGAAATGGATGATGGCAATGAAATATCCAATCACAGCAGAC<br>AAGAGAATAATGGAAATGATTCCTGAAAGAAATGAACAAGGCCA<br>GACGCTTTGGAGCAAGACAAATGATGCTGGATCAGACAGAGTGAT<br>GGTGTCTCCCCTAGCTGTAACTTGGTGGAATAGAAATGGACCGAC<br>AGCAAGTACAGTCCATTATCCAAAGGTCTACAAAACATACTTTGA<br>GAAGGTTGAAAGGTTAAAGCATGGAACCTTCGGTCCCGTTCACTTC<br>CGAAACCAAATTAAAATACGCCGCCGAGTTGACATAAACCCAGGC<br>CACGCAGATCTCAGTGCCAAAGAAGCACAAGATGTCATCATGGAG<br>GTCGTTTTCCCAAATGAAGTGGGAGCTAGAATATTGACATCAGAG<br>TCACAATTGACAATAACGAAAGAGAAAAAGAAGAACTCCAGGA<br>TTGCAAGATTGCTCCTTTAATGGTGGCATACATGTTGGAAAGAGAA<br>CTGGTCCGCAAAACCAGATTCCTACCAGTAGCAGGTGGGACAAGC<br>AGTGTGTACATTGAGGTACTGCACCTGACCCAAGGGACCTGCTGG<br>GAACAGATGTACACTCCAGGCGGAGAAGTGAGAAATGACGATGTT<br>GACCAGAGTTTGATCATCGCGGCCAGAAACATTGTTAGGAGAGCA<br>ACGGTATCAGCGGATCCACTGGCATCATTATTGGAGATGTGCCAC<br>AGCACACAAATTGGTGGGACAAGGATGGTGGATATCCTTAGGCAA<br>AATCCAACTGAGGAACAAGCTGTGGATATATGCAAAGCAGCAATG<br>GGTTTAAGGATTAGTTCATCCTTTAGCTTTGGAGGATTCACCTTCA<br>AAAGAACAAGTGGTTCATCCATTAGAAAGGAAGAGGAAGTGCTTA<br>CAGGCAACCTCCAAACATTGAAAATAAGAGTACATGAGGGGTAT<br>GAGGAGTTCACAATGGTTGGGCGAAGAGCAACAGCCATTCTAAGG<br>AAAGCAACTAGAAGGCTGATTCAGTTGATAGTAAGTGGAAGAGAC<br>GAACAATCAATCGCTGAAGCAATCATCGTAGCCATGGTGTTCTCAC<br>AGGAGGATTGCATGATAAAGGCAGTCCGAGGCGATCTAAATTTT<br>GTGAACAGAGCAAACCAAAGATTGAACCCCATGCATCAACTCCTG<br>AGACACTTCCAAAAAGATGCAAAAGTGCTGTTTCAAAATTGGGGGG<br>ATCGAACCCATTGATAATGTCATGGGGATGATTGGAATATTGCCTG<br>ACATGACTCCAAGCACAGAGATGTCACTAAGAGGAGTAAGAGTT<br>AGTAAAATGGGAGTAGATGAATATTCCAGCACTGAGAGAGTGGTT<br>GTAAGCATTGACCGTTTCTTGCGGGTTCGAGATCAGCAGGGGAAC<br>GTACTCCTATCTCCCGAAGAAGTCAGCGAAACACTGGGAACAGAA<br>AAATTAACAATAACATATTCATCATCAATGATGTGGGAAATCAAT |

| | | | |
|---|---|---|---|
| | | | GGTCCTGAGTCAGTGCTGGTCAACACCTATCAATGGATCATCAGA<br>AATTGGGAGATTGTGAAGATTCAATGGTCTCAAGACCCCACGATG<br>CTGTACAATAAGGTGGAGTTTGAACCGTTCCAATCCTTGGTACCTA<br>AAGCTGCCAGAGGCCAATACAGTGGATTTGTGAGAACACTGTTC<br>CAACAAATGCGTGACGTATTGGGGACATTTGATACTATTCAGATA<br>ATAAAGCTGTTACCGTTTGCAGCAGCCCCACCGGAGCATAGCAGA<br>ATGCAATTTTCTTCCCTGACCGTGAACGTAAGAGGCTCGGGAATGA<br>GAATACTCGTAAGGGGTAACTCCCCTGTGTTCAACTACAATAAG<br>GCAACCAAAAGGCTTGCCGTCCTTGGAAAGGACGCAGGTGCATTA<br>ACAGAGGATCCAGATGAGGGGACAACAGGAGTGGAATCTGCAGT<br>GCTGAGGGGGTTCCTAATTCTGGGCAGGGAGGACAGAAGATATGG<br>ACCAGCACTAAGCATCAATGAACTGAGCAATCTTGCGAAAGGGGA<br>GAAAGCCAATGTGCTGATAGGGCAAGGAGACGTGGTGCTGGTAAT<br>GAAACGGAAACGGGACTCTAGCATACTTACTGACAGCCAGACAGC<br>GACCAAAAGAATTCGGATGGTCATCAATTAGTATCGAG |
| DETU | PB1 | 1-2277 | **ATG**GATGTCAACCCGACTTTACTCTTCTTGAAAGTGCCAGCGCAAA<br>ATGCTATAAGTACCACATTCCCTTATACTGGAGATCCTCCATACAG<br>CCATGGAACAGGAACAGGATACACCATGGACACAGTCAACAGAA<br>CGCATCAATACTCAGAAAAGGGAAAGTGGACAAAAAACACCGAG<br>ACTGGAGCACCCCAACTCAACCCAATTGATGGACCATTACCTGAG<br>GATAACGAGCCAAGCGGATATGCACAAACGGATTGTGTGTTGGAA<br>GCAATGGCTTTCCTTGAAGAGTCCCACCCAGGGATCTTTGAAAACT<br>CATGTCTTGAAACAATGGAAATTGTTCAACAAACAAGAGTGGAC<br>AAACTGACCCAAGGTCGTCAGACCTATGACTGGACATTGAATAGA<br>AACCAGCCGGCTGCAACTGCTTTAGCCAACACTATAGAAGTCTTCA<br>GATCGAACGGTCTAACAGCCAATGAGTCAGGGAGACTGATAGATT<br>TCCTCAAAGATGTGATGGAGTCAATGGACAAAGAAGAAATGGAA<br>ATAACAACACATTTCCAAAGAAAGAGAAGAGTAAGAGACAATAT<br>GACCAAGAAAATGGTCACACAAAGAACAATAGGGAAGAAAAAAC<br>AGAGACTGAACAAGAAGAACTACTTGGTAAGGGCACTGACACTGA<br>ACACAATGACAAAAGATGCAGAAAGAGGCAAGTTGAAGAGGCGG<br>GCAATTGCAACACCCGGGATGCAAATCAGAGGGTTCGTGTACTTT<br>GTCGAAACATTAGCGAGGAGCATCTGCGAGAAACTTGAGCAATCT<br>GGGCTCCCTGTTGGAGGAAATGAAAAAAAGGCTAAGTTGGCAAAT<br>GTCGTGAGAAAGATGATGACTAACTCACAAGACACAGAGCTATCC<br>TTTACAATTACTGGAGACAATACCAAGTGGAACGAGAATCAGAAT<br>CCTCGGATTTTTTTGGCAATGATAACATATATCACAAGAAATCAAC<br>CTGAGTGGTTTAGAAATGTGTTAAGTATTGCCCCTATAATGTTCTC<br>AAACAAAATGGCAAGATTAGGGAAAGGATACATGTTCGAAAGTA<br>AGAGCATGAAGCTACGGACACAAATACCAGCAGAAATGCTTGCAA<br>CCATTGACCTGAAATATTTCAACGAATCGACAAGAAAGAAAATTG<br>AGAAAATAAGGCCTCTCCTAATAGAAGGGACAGCCTCGTTGAGTC<br>CTGGAATGATGATGGGCATGTTCAACATGCTGAGTACAGTCTTGG<br>GAGTATCAATTCTAAATCTTGGCCAAAAGAGGTACACCAAAACCA<br>CATACTGGTGGGACGGACTCCAATCCTCTGATGATTTCGCTCTCAT<br>AGTAAATGCACCGAATCATGAGGGAATACAGGCAGGAGTGGACA<br>GGTTCTATAGGACTTGTAAATTGGTTGGGATCAATATGAGTAAAA<br>AGAAATCCTATATAAATCGGACAGGAACATTTGAATTCACAAGCT<br>TTTTCTACCGTTATGGGTTTGTAGCCAACTTCAGCATGGAGCTGCC<br>CAGCTTTGGAGTTTCTGGGATTAATGAATCGGCTGACATGAGCATT<br>GGAGTTACAGTAATAAAGAATAACATGATAAACAACGATCTTGGA<br>CCAGCAACAGCTCAAATGGCTCTTCAGCTATTTATCAAGGACTACA<br>GATATACATATCGATGCCACAGGGGTGATACACAAATACAAACAA<br>GGAGATCATTCGAGCTAAAGAAGCTGTGGGAGCAGACCCGTTCAA<br>AGGCAGGACTGTTGGTTTCAGATGGAGGCCCAAACTTATACAAT<br>ATACGGAATCTCCACATCCCAGAGGTCTGCTTGAAGTGGGAACTG |

| | | | |
|---|---|---|---|
| | | | ATGGATGAAGATTACCAGGGTAGACTTTGTAATCCCCTGAACCCCT TTGTCAGTCATAAGGAAATTGAATCCGTAAACAATGCTGTAGTGAT GCCAGCCCATGGTCCGGCCAAAAGCATGGAATATGATGCTGTT GCGACCACACACTCATGGGTCCCTAAGAGGAACCGTTCCATTCTG AATACCAGTCAAAGAGGAATCCTTGAGGATGAACAGATGTATCAG AAGTGCTGCAATCTATTTGAAAAATTCTTCCCTAGTAGCTCATACA GGAGGCCAGTTGGAATCTCCAGTATGGTGGAGGCCATGGTGTCT AGGGCCCGAATTGATGCACGGATTGACTTCGAGTCTGGTAGGATT AAGAAGGAAGAGTTTGCTGAGATCATGAAGATCTGTTCCACCATT GAAGAGATCAGACGGCAAAAACAGTGA |
| DETU | PA | 7-2190 | GACTTTGTGCGACAATGCTTCAATCCAATGATCGTCGAGCTTGCGG AAAAGACAATGAAAGAATATGGGGAAAATCCAAAAATCGAAACG AACAAATTCGCTGCAATATGCACTCACTTAGAGGTCTGTTTCATGT ATTCGGATTTCCACTTTATTGATGAACGAGGTAAATCAATAATTGT AGAATCTGGCGATCCGAATGCATTATTGAAACACCGATTTGAGAT AATTGAAGGGAGGGACCGAACGATGGCTTGGACAGTGGTAAATAG TATCTGCAACACCACAGGAGTCGATAAGCCTAAATTCCTCCCAGAT TTGTATGATTACAAGGAGAACCGATTCATTGAAATTGGAGTGACA AGGAGGGAAGTTCACACATACTACCTAGAAAAGGCAAATAAGATA AAATCAGAGAAGACACACATTCACATATTCTCATTCACTGGGGAG GAGATGGCCACCAAAGCTGATTATATCCTTGATGAAGAGAGCAGA GCAAGGATCAAAACCAGGTTGTTCACTATCAGGCAAGAAATGGCC AATAGGGGTCTGTGGGATTCCTTTCGTCAATCTGAGAGAGGCGAA GAGACAATTGAAGAAAGGTTTGAAATCACAGGAACCATGCGCAGG CTTGCCGACCAAAGTCTCCCACCGAATTTCTCCAGCCTTGAAAATT TTAGAGCCTATGTGGATGGATTCAAACCGAACGGCTGCCTTGAGG GCAAGCTTTCTCAAATGTCAAAAGAAGTGAACGCCAGAATTGAGC CATTCATGAAGACAACACCACGCCCTCTCAGATTACCTGATGGTCC TCCTTGCTCTCAGCGGTCGAAATTCTTACTGATGGATGCTCTTAAA TTGAGCATCGAAGACCCAAGCCATGAGGGAGAAGGTATACCGCTA TATGATGCAATCAAATGCATGAAGACGTTTTTTGGTTGGAAAGAG CCCAACATTGTAAAACCACATGTAAAAGGCATAAATCCCAACTAT CTCTTGGCTTGGAAGCAGGTGCTGGTAGAACTCCAAGACATTGAA AATGAAGAGAAATCCCAAAAACAAAAAACATGAAGAAAACAAG CCAACTAAAGTGGGCACTCGGTGAGAATATGGCACCTGAAAAAGT GGACTTTGAGGACTGCAAAGATGTTAGCGATCTAAGACAGTATGA CAGTGATGAACCAGAGCCCAGATCATTATCAAGCTGGATCCAGAG CGAATTCAACAAAGCATGCGAATTGACAGATTCGAGTTGGATTGA ACTTGATGAAATAGGAGAAGATGTTGCTCCAATTGAGCACATTGC GAGTATGAGAAGAAACTACTTCACAGCGGAAGTGTCTCATTGCAG GGCTACTGAATATATAATGAAAGGAGTTTATATAAATACAGCCCT GTTGAATTCATCCTGTGCAGCCATGGATGACTTCCAATTGATTCCA ATGATAAGCAAGTGCAGAACCAAAGAAGGAAGACGGAAGACAAA TCTATATGGGTTCATTATAAAAGGAAGATCCCATTTGAGGAATGAT ACCGATGTGGTAAATTTTGTGAGCATGGAGTTCTCTCTTACTGACC CGAGGCTGGAACCACACAAGTGGGAAAAGTACTGTGTTCTCGAAA TAGGAGACATGCTCCTACGAACTGCAATAGGCCAAGTATCAAGAC CCATGTTTCTTTATGTAAGGACCAATGGGACTTCCAAGATCAAGAT GAAATGGGGCATGGAGATGAGGCGATGCCTTCTTCAATCCCTCCA ACAAATTGAGAGCATGATTGAGGCAGAATCTTCTGTCAAAGAGAA GGACATGACCAAGGAATTCTTTGAAAATAAATCAGAAACGTGGCC AATTGGGGAATCACCTAAGGGGGTGGAGGAAAGCTCTATTGGGAA AGTGTGTAGAACATTACTAGCAAAATCTGTATTCAACAGCCTATAT GCATCTCCACAACTTGAGGGGTTTTCAGCTGAGTCGAGAAAGTTAC TTCTCATTGTTCAGGCATTTAGGGACAACCTGGAACCTGGGACCTT CGATCTTGGGGGGGCTATATGAAGCAATTGAGGAGTGCCTGATTAA |

| | | | TGATCCCTGGGTTTTGCTTAATGCATCTTGGTTCAACTCCTTCCTTA CACATGCACTGAAATAGTTGTGGCAATGCTACTATTTGCTATCCAT ACTGTCCAA |
|---|---|---|---|
| DETU | HA | 1-1728 | **ATG**GAGAAAATAGTGCTTCTTCTTGCAGTGGTTAGCCTTGTTAAAA GTGATCAGATTTGCATTGGTTACCATGCAAACAACTCAACAAAAC AGGTTGACACAATAATGGAAAAAAACGTCACTGTTACACATGCCC AAGACATACTGGAAAAGACACACAACGGGAAGCTCTGCGATCTTA ATGGAGTGAAGCCCCTGATTCTAAAGGATTGTAGCGTAGCTGGGT GGCTCCTTGGAAATCCAATGTGCGACGAGTTTATCAGGGTGCCGG AATGGTCTTACATCGTGGAGAGGGCTAACCCAGCCAACGACCTCT GTTACCCAGGGACCCTCAATGACTATGAGGAACTGAAACACCTA CTGAGCAGAATAAATCATTTTGAGAAAACTCTGATCATCCCCAAG AGTTCTTGGCCCAATCATGAAACATCATTAGGGGTGAGCGCAGCA TGTCCATACCAGGGAGCATCCTCATTTTTCAGAAATGTGGTATGGC TCATCAAAAAGAACGATGCATACCCGACAATAAAGATAAGCTAC AATAATACCAATCGGGAAGATCTTTTGATACTGTGGGGGATTCATC ATCCCAACAATGCAGAAGAGCAGACAAATCTCTATAAAAACCCAG ACACTTATGTTTCCGTTGGGACATCAACATTAAACCAGAGATTGGT GCCAAAAATAGCTACTAGATCCCAAGTAAACGGGCAACGTGGAAG AATGGATTTCTTCTGGACAATTTTAAAACCGAATGATGCAATCCAC TTTGAGAGTAATGGAAATTTCATTGCTCCAGAATATGCCTACAAAA TTGTCAAGAAAGGGGACTCAACAATTATGAAAAGTGAAGTGGAGT ATGGCCACTGCAACACCAAATGTCAAACCCCAATAGGGGCGATAA ACTCTAGCATGCCATTCCACAATATACACCCTCTCACCATCGGGGA ATGCCCCAAATATGTGAAGTCAAACAAATTAGTCCTTGCGACTGG GCTCAGAAATAGTCCTCTAAGGGAAAGAAGAAGAAAAAGAGGAC TATTTGGAGCTATAGCAGGGTTTATAGAGGGAGGATGGCAGGGAA TGGTAGACGGTTGGTATGGGTACCACCATAGCAATGAGCAGGGGA GTGGGTACGCTGCAGACAAAGAATCCACCCAAAAGGCAGTAGATG GAGTTACCAATAAGGTCAACTCAATCATTGACAAAATGAACACTC AATTTGAGGCCGTTGGAAGGGAATTTAATAACTTAGAAAGGAGAA TAGAGAATTTAAACAAGAAAATGGAAGACGGATTCCTAGATGTCT GGACTTATAATGCTGAACTTTTAGTTCTTATGGAAAATGAGAGAAC TCTAGATTTCCATGACTCAAATGTCAAGAACCTTTACGACAAAGTC CGACTACAGCTTAGGGATAATGCAAAGGAGCTGGGTAATGGTTGT TTCGAGTTCTATCACAAATGTGATAACGAATGTATGGAAAGCGTA AGAAATGGGACGTATGACTACCCTAAGTATTCAGAAGAAGCAAGA TTAAAAAGAGAAGAAATAAGCGGAGTGAAATTAGAATCAATAGG AACTTACCAAATACTGTCAATTTATTCAACAGTGGCGAGTTCCCTA GCACTGGCAATCATAGTGGCTGGTCTATCTTTATGGATGTGCTCTA ATGGGTCGCTACAATGCAGAATTTGCATCTAAATTTGTGAGCTCAG ATTGTAATTA |
| DETU | NP | 1-1497 | **ATG**GCGTCTCAAGGCACCAAACGATCTTATGAACAGATGGAAACT GGTGGAGAACGCCAGAATGCCACTGAAATCAGAGCATCTGTTGGA AGAATGGTTGGTGGAATTGGAAGGTTTTATATACAGATGTGCACT GAACTCAAACTCAGCAATTATGAGGGGAGACTGATCCAGAACAGC ATAACAATAGAAAGAATGGTTCTCTCTGCATTTGATGAAAGGAGG AACAAGTACCTGGAAGAACATCCCAGTGCGGGGAAGGACCCAAA GAAAACTGGAGGTCCAATCTACAGAAGAAGAGACGGAAAGTGGA TGAGGGAGCTGATTCTGTATGACAAAGAAGAGATCAGAAGGATCT GGCGTCAAGCAAATAATGGAGAAGATGCAACTGCTGGTCTCACCC ATCTGATGATCTGGCACTCCAACCTGAATGATGCCACATATCAGAG GACAAGGGCTCTCGTGCGCACTGGAATGGATCCCAGAATGTGCTC TCTGATGCAAGGATCAACTCTCCCAAGAAGGTCTGGAGCTGCTGG TGCAGCAGTAAAAGGGGTCGGAACAATGGTAATGGAATTGATTCG AATGATAAAGCGAGGGATTAATGATCGGAATTTCTGGAGAGGCGA |

| | | | |
|---|---|---|---|
| | | | AAATGGAAGAAGGACAAGGATTGCCTATGAGAGAATGTGCAACAT CCTCAAAGGGAAATTTCAAACAGCAGCACAAAGAGCAATGATGGA TCAAGTGCGAGAAAGCAGGAATCCTGGGAATGCTGAAATTGAAGA TCTCATTTTTCTGGCACGGTCTGCACTCATCCTGAGAGGATCAGTG GCCCACAAGTCTTGTCTGCCTGCTTGTGTTTACGGACTTGCTGTGG CCAGTGGATATGACTTTGAGAGAGAAGGATACTCTCTGGTTGGA ATAGACCCTTTCCGTCTGCTTCAAAACAGCCAGGTCTTCAGTCTCA TTAGACCAAATGAAAACCCAGCACATAAAAGTCAGTTGGTATGGA TGGCATGCCATTCAGCAGCGTTTGAGGACCTGAGGGTATCAAGTTT CATCAGAGGGACAAGAGTGGTCCCAAGAGGACAACTATCCACCAG AGGAGTTCAAATTGCATCAAATGAAAACATGGAAACAATGGACTC CAGCACTCTTGAATTGAGGAGCAGATACTGGGCTATAAGAACCAG GAGTGGAGGAAACACCAACCAACAGAGAGCTTCTGCAGGACAAA TCAGCGTACAACCCACCTTCTCAGTACAGAGAAATCTTCCCTTTGA AAGAGCGACCATCATGGCGGCATTTACAGGGAACACTGAAGGCAG GACCTCTGACATGAGGACTGAGATCATAAGAATGATGGAAAGTGC CAAACCAGAAGATGTGTCCTTCCAGGGGCGGGGAGTCTTCGAGCT CTCGGACGAAAAGGCAACGAACCCGATCGTGCCTTCCTTTGACAT GAGCAACGAAGGATCTTATTTCTTCGGAGACAGTGCAGAGGAGTA TGACAATTAA |
| DETU | NA | 1-1413 | **ATG**AATCCAAATCAGAAAATAGTAACCATTGGCTCCATTTCATTA GGGTTGGTTGTATTCAATGTTCTACTGCATGCTGTGAGCATCATAT TAACAGTGTTAGCCCTGGGGAAGAGTGAAAACAATGGAATCTGCA ATGGAACTGTAGTGAGGGAATACAATGAAACAGTTAGAATAGAG AAAGTGACTCAATGGTACAATACTAGCGTAGTCGAATATGTACCG CATTGGAATGAGGGCACTTATATAAATAACACCGAACCAATATGT GATGTCAAGGGCTTTGCACCTTTTTCCAAGGACAACGGGATAAGA GTTGGCTCCAGGGGACATATTTTTGTCATAAGAGAGCCTTTCGTC TCTTGTTCACCTGTAGAGTGCAGGACTTTCTTCCTCACTCAGGGAT CTCTACTCAATGACAAACACTCAAATGGAACAGTGAAGGATAGAA GCCCATTCAGAACTCTCATGAGTGTCGAAGTGGGCCAATCACCCA ATGTATATCAAGCCAGGTTTGAAGCTGTGGCATGGTCAGCAACA GCCTGTCATGATGGTAAGAAGTGGATGACGATTGGTGTAACAGGG CCAGATTCTAAAGCAGTAGCAGTAGTTCATTACGGAGGGGTGCCT ACTGACGTTGTTAACTCCTGGGCAGGAGATATATTAAGAACTCAG GAGTCATCTTGTACTTGCATTCAAGGTAATTGTTATTGGGTAATG ACTGACGGTCCTGCCAATAGACAGGCGCAGTATAGAATATACAAA GCAAATCAAGGCAAAATAATTGGCCGAACAGATGTTAGCTTTAGT GGAGGACATATTGAGGAATGTTCTTGTTATCCAAATGATGGTAAA GTGGAATGCGTGTGTAGAGACAACTGGACGGGAACTAACAGGCCT GTACTAATTATTTCGCCTGATCTCTCTTACAGGGGTTGGGTATTTATG TGCAGGGTTGCCCAGTGACACTCCAAGAGGGGAAGATACTCAATT TGTCGGTTCATGCACTAGTCCCATGGGAAATCAGGGATATGGCGT AAAAGGGTTCGGGTTTCGACAGGGAACTGATGTGTGGGTGGGG CGGACAATTAGTCGAACCTCCAGATCAGGATTTGAAATAATAAGG ATAAAGAATGGTTGGACGCAAACAAGCAAAGAACAGATTAGAAG ACAGGTGGTTGTTGATAACTCGAATTGGTCGGGATACAGTGGGTCT TTCACTTTACCAGTAGAATTGTCTGGGAGGGAATGTTTGGTTCCC TGTTTTTGGGTCGAAATGATCAGAGGTAGGCCAGAAGAGAGAACA ATCTGGACCTCTAGTAGCTCCATTGTAATGTGTGGAGTTGATTATG AAATTGCCGATTGGTCATGGCACGATGGAGCTATTCTTCCCTTTGA CATCGATAAGACGTAA |
| DETU | MP | -1-982 | **G**ATGAGTCTTCTAACCGAGGTCGAAACGTACGTTCTCTCTATCATC CCGTCAGGCCCCCTCAAAGCCGAGATCGCGCAGAGACTTGAAGAT GTCTTTGCAGGGAAAAACACCGATCTCGAGGCTCTCATGGAGTGG CTAAAGACAAGACCAATCCTGTCACCTCTGACTAAAGGGATTTTG |

| | | | |
|---|---|---|---|
| | | | GGATTTGTGTTCACGCTCACCGTGCCCAGTGAGCGAGGACTGCAG<br>CGTAGACGCTTCGTCCAGAATGCCCTAAATGGAAACGGGGATCCA<br>AATAATATGGATAAGGCAGTTAAGCTATATAAGAAGCTGAAAAGA<br>GAGATAACATTCCATGGGGCTAAGGAGGTCGCACTTAGCTACTCA<br>ACCGGTGCACTTGCCAGCTGCATGGGTCTCATATACAACAGGATG<br>GGAACGGTGACTACAGAAGTGGCTTTTGGCCTAGTGTGTGCCACTT<br>GTGAGCAGATTGCAGATTCACAGCATCGGTCCCACAGACAGATGG<br>CAACCATCACCAACCCATTAATCAGACATGAGAACAGAATGGTGC<br>TGGCCAGCACTACAGCTAAGGCCATGGAGCAGATGGCAGGATCAA<br>GCGAGCAGGCATCAGAAGCCATGGAGGTTGCTAATCAGGCCAGGC<br>AGATGGTACAGGCAATGAGGACAATTGGGACTCATCCTAACTCTA<br>GTGCTGGTCTGAGAGATAATCTTCTTGAAAATTTGCAGGCCTACCA<br>GAACCGAATGGGAGTGCAGATGCAGCGATTCAAGTGATCCTCTTG<br>TTGTTGCCGCAAATATCATTGGGATCCTGCACTTGATATTGTGGAT<br>CCTTGATCGTCTTTTCTTCAAATGCATTTATCGTCGCCTTAAATACG<br>GTTTGAAAATAGGGCCTTCTACGGAAGGGGTACCTGAGTCTATGA<br>GGGAAGAGTACCGGCAGGAACAGCAGAGTGCTGTGGATGTTGACG<br>ATGGTCATTTTGTCAACATAGAATTGGAGTAA |
| DETU | NS | 2-838 | **TG**GACTCCAACACTGTGTCAAGCTTTCAGGTAGACTGCTTTCTTTG<br>GCATGTCCGCAAACGATTTGCAGACCAAGAACTGGGTGATGCCCC<br>ATTCCTTGACCGGCTTCGCCGAGACCAGAAGTCCCTAAGAGGAAG<br>AGGCAGCACTCTTGGTCTGGACATCGAGACAGCTACTCGTGCGGG<br>AAAGCAAATATTGGAGCGGATTCTGGGGGAAGAATCTGATGAAGC<br>ACTTAAAATGAATATTGCTTCTGTACCGACTTCACGCTACCTAACT<br>GACATGACTCTTGAAGAAATGTCAAGAGACTGGTTCATGCTCATG<br>CCCAAGCAGAAAGTAGCAGGTTCTCTCTGCATCAAAATGGACCAG<br>GCAATAATGGATAAAACCATCATACTGAAAGCAAACTTCAGTGTG<br>ATTTTTGATCGGCTGGAAACCCTAATATTACTTAGAGCTTTCACAG<br>AAGAAGGAGCAATTGTGGGAGAAATCTCACCATTACCTTCTCTTCC<br>AGGACATACTGATGAGGATGTCAAAATTGCAATTGGGGTCCTCAT<br>CGGAGGGCTTGAATGGAATGATAACACAGTTCGAGTCTCTGAAAC<br>TCTACAGAGATTCACTTGGAGAAGCAGTAATGAGGATGGGAGACC<br>TTCACTCCCTTCAAAACAGAAACGGAAAATGGCGAGAACAATTGA<br>GTCAGAAGTTCGAGGAAATAAGATGGCTGATTGAGGAAATGCGAC<br>ATAGATTGAAGATCACAGAGAACAGCTTCGAACAAATAACGTTTA<br>TGCAAGCTTTACAACTATTGCTTGAAGTGGAGCAAGAGATAAGAA<br>CCTTCTCGTTTCAGCTTATTTAA |
| UKDD | PB2 | 1-2298 | **ATG**GAGAGAATAAAAGAACTAAGAGATCTAATGTCTCAATCCCGC<br>ACTCGCGAGATACTAACAAAAACCACTGTGGACCATATGGCCATA<br>ATCAAGAAATACACATCAGGAAGACAAGAGAAGAACCCTGCTCTC<br>AGAATGAAATGGATGATGGCAATGAAATATCCAATCACAGCAGAC<br>AAGAGAATAATGGAAATGATTCCTGAAAGAAATGAACAAGGCCA<br>GACGCTTTGGAGTAAGACAAATGATGCTGGATCAGACAGAGTGAT<br>GGTGTCTCCCCTAGCTGTAACTTGGTGGAATAGAAATGGACCGAC<br>AGCAAGTACAGTCCATTATCCAAAGGTCTACAAAACATACTTTGA<br>GAAGGTTGAAAGGTTAAAGCATGGAACCTTCGGTCCCGTTCACTTC<br>CGAAACCAAATTAAAATACGCCGCCGAGTTGACATAAACCCAGGC<br>CACGCAGATCTCAGTGCCAAAGAAGCACAAGATGTCATCATGGAG<br>GTCGTTTTCCCAAATGAAGTGGGAGCTAGAATATTGACATCAGAG<br>TCACAATTGACAATAACGAAAGAGAAAAAGAAGAACTCCAGGA<br>TTGCAAGATTGCTCCTTTAATGGTGGCATACATGTTGGAAAGAGAA<br>CTGGTCCGCAAAACCAGATTCCTACCAGTAGCAGGTGGGACAAGC<br>AGTGTGTACATTGAGGTACTGCACTTGACCCAAGGGACCTGCTGG<br>GAACAGATGTACACTCCAGGCGGAGAAGTGAGAAATGACGATGTT<br>GACCAGAGTTTGATCATCGCGGCCAGAAACATTGTTAGGAGAGCA<br>ACGGTATCAGCGGATCCACTGGCATCATTATTGGAGATGTGCCAC |

| | | | |
|---|---|---|---|
| | | | AGCACACAAATTGGTGGGACAAGGATGGTGGATATCCTTAGGCAA AATCCAACTGAGGAACAAGCTGTGGATATATGCAAAGCAGCAATG GGTTTAAGGATTAGTTCATCCTTTAGCTTTGGAGGATTCACCTTCA AAAGAACAAGTGGTTCATCCATTAGAAAGGAAGAGGAAGTGCTTA CAGGCAACCTCCAAACATTGAAAATAAGAGTACATGAGGGGTATG AGGAGTTCACAATGGTTGGGCGAAGAGCAACAGCCATTCTAAGGA AAGCAACTAGAAGGCTGATTCAGTTGATAGTAAGTGGAAGAGACG AACAATCAATCGCTGAAGCAATCATCGTAGCCATGGTGTTCTCACA GGAGGATTGCATGATAAAGGCAGTCCGAGGCGATCTAAATTTTGT GAACAGAGCAAACCAAAGATTGAACCCCATGCATCAACTCCTGAG ACACTTCCAAAAAGATGCAAAAGTGCTGTTTCAAAATTGGGGGGAT TGAACCCATTGATAATGTCATGGGGATGATTGGAATATTGCCTGAC ATGACTCCAAGCACAGAGATGTCACTAAGAGGAGTAAGAGTTAGT AAAATGGGAGTAGATGAATATTCCAGCACTGAGAGAGTGGTTGTA AGCATTGACCGTTTCTTGCGGGTTCGAGATCAGCAGGGGAACTTAC TCCTATCTCCCGAAGAAGTCAGCGAAACACTGGGAACAGAAAAGT TAACAATAACATATTCATCATCAATGATGTGGGAAATCAATGGTCC TGAGTCAGTGCTGGTCAACACCTATCAATGGATCATCAGAAATTG GGAGATTGTGAAGATTCAATGGTCTCAAGACCCCACGATGCTGTA CAATAAGGTGGAGTTTGAACCGTTCCAATCCTTGGTACCTAAAGCT GCCAGAGGCCAATACAGTGGATTTGTGAGAACACTGTTCCAACAA ATGCGTGACGTATTGGGGACATTTGATACTATTCAGATAATAAAGC TGTTACCGTTTGCAGCAGCCCCACCGGAGCATAGCAGAATGCAAT TTTCTTCCCTGACCGTGAATGTAAGAGGCTCGGGAATGAGAATACT CGTAAGGGGTAACTCCCCTGTGTTCAACTACAATAAGGCAACCAA AAGGCTTGCCGTCCTTGGAAAGGACGCAGGTGCATTAACAGAGGA TCCAGATGAGGGGACAACAGGAGTGGAATCTGCAGTGCTGAGGGG GTTCCTAATTCTGGGCAGGGAGGACAGAAGATATGGACCAGCACT AAGCATCAATGAACTGAGCAATCTTGCGAAAGGGGAGAAAGCCA ATGTGCTGATAGGGCAAGGAGACGTGGTGCTGGTAATGAAACGGA AACGGGACTCTAGCATACTTACTGACAGCCAGACAGCGACCAAAA GAATTCGGATGGTCATCAATTAGTATCGAGTTGTTTAAAAA |
| UKDD | PB1 | 1-2277 | **ATG**GATGTCAACCCGACTTTACTCTTCTTGAAAGTGCCAGCGCAAA ATGCTATAAGTACCACATTCCCTTATACTGGAGATCCTCCATACAG CCATGGAACAGGAACAGGATACACCATGGACACAGTCAACAGAA CGCATCAATACTCAGAAAAGGGAAAGTGGACAAAAAACACCGAG ACTGGAGCACCCCAACTCAACCCAATTGATGGACCATTACCTGAG GATAACGAGCCAAGCGGATATGCACAAACGGATTGTGTGTTGGAA GCAATGGCTTTCCTTGAAGAGTCCCACCCAGGGATCTTTGAAAACT CATGTCTTGAAACAATGGAAATTGTTCAACAAACAAGAGTGGACA AACTGACCCAAGGTCGTCAGACCTATGACTGGACATTGAATAGAA ACCAGCCGGCTGCAACTGCTTTAGCCAACACTATAGAAGTCTTCAG ATCGAACGGTCTAACAGCCAATGAGTCAGGGAGACTGATAGATTT CCTCAAAGATGTGATGGAGTCAATGGACAAAGAAGAAATGGAAAT AACAACACATTTCCAAAGAAAGAGAAGAGTAAGAGACAATATGA CCAAGAAAATGGTCACACAAGAACAATAGGGAAGAAAAAACAG AGACTGAACAAGAAGAACTACTTGGTAAGGGCACTGACACTGAAC ACAATGACAAAAGATGCAGAAAGAGGCAAGTTGAAGAGGCGGGC AATTGCAACACCCGGGATGCAAATCAGAGGGTTCGTGTACTTTGTC GAAACATTAGCGAGGAGCATCTGCGAGAAACTTGAGCAATCTGGG CTCCCTGTTGGAGGAAATGAAAAAAAGGCTAAGTTGGCAAATGTC GTGAGAAAGATGATGACTAACTCACAAGACACAGAGCTATCCTTT ACAATTACTGGAGACAATACCAAGTGGAACGAGAATCAGAATCCT CGGATTTTTTTGGCAATGATAACATATATCACAAGAAATCAACCTG AGTGGTTTAGAAATGTGTTAAGTATTGCCCCTATAATGTTCTCAAA CAAAATGGCAAGATTAGGGAAAGGATACATGTTCGAAAGTAAG |

| | | | AGCATGAAGCTACGGACACAAATACCAGCAGAAATGCTTGCAACC<br>ATTGACCTGAAATATTTCAACGAATCGACAAGAAAGAAAATTGAG<br>AAAATAAGGCCTCTCCTAATAGAAGGAACAGCCTCGTTGAGTCCT<br>GGAATGATGATGGGCATGTTCAACATGCTGAGTACAGTCTTGGGA<br>GTATCAATTCTAAATCTTGGCCAAAAGAGGTACACCAAAACCACA<br>TACTGGTGGGACGGACTCCAATCCTCTGATGATTTCGCTCTCATAG<br>TAAATGCACCGAATCATGAGGGAATACAGGCAGGAGTGGACAGGT<br>TCTATAGGACTTGTAAATTGGTTGGGATCAATATGAGTAAAAAG<br>AAATCCTATATAAATCGGACAGGAACATTTGAATTCACAAGCTTTT<br>TCTACCGTTATGGGTTTGTAGCCAACTTCAGCATGGAGCTGCCCAG<br>CTTTGGAGTTTCTGGGATTAATGAATCGGCTGACATGAGCATTGGA<br>GTTACAGTAATAAAGAATAACATGATAAACAACGATCTTGGACCA<br>GCAACAGCTCAAATGGCTCTTCAGCTATTTATCAAGGACTACAGAT<br>ATACATATCGATGCCACAGGGGTGATACACAAATACAAACAAGGA<br>GATCATTCGAGCTAAAGAAGCTGTGGGAGCAGACCCGTTCAAAGG<br>CAGGACTGTTGGTTTCAGATGGAGGCCCAAACTTATACAATATAC<br>GGAATCTCCACATCCCAGAGGTCTGCTTGAAGTGGGAACTGATGG<br>ATGAAGATTACCAGGGTAGACTTTGTAATCCCCTGAACCCCTTTGT<br>CAGTCATAAGGAAATTGAATCCGTAAACAATGCTGTAGTGATGCC<br>AGCCCATGGTCCGGCCAAAAGCATGGAATATGATGCTGTTGCGAC<br>CACACACTCATGGGTCCCTAAGAGGAACCGTTCCATTCTGAATACC<br>AGTCAAAGAGGAATCCTTGAGGATGAACAGATGTATCAGAAGTGC<br>TGCAATCTATTTGAAAAATTCTTCCCTAGTAGCTCATACAGGAGGC<br>CAGTTGGAATCTCCAGTATGGTGGAGGCCATGGTGTCTAGGGCCC<br>GAATTGATGCACGGATTGACTTCGAGTCTGGTAGGATTAAGAAGG<br>AAGAGTTTGCTGAGATCATGAAGATCTGTTCCACCATTGAAGAGA<br>TCAGACGGCAAAAACAGTGA |
| UKDD | PA | 1-2151 | **ATG**GAAGACTTTGTGCGACAATGCTTCAATCCAATGATCGTCGAG<br>CTTGCGGAAAAGACAATGAAAGAATATGGGGAAAATCCAAAAAT<br>CGAAACGAACAAATTCGCTGCAATATGCACTCACTTAGAGGTCTG<br>TTTCATGTATTCGGATTTCCACTTTATTGATGAACGAGGTAAATCA<br>ATAATTGTAGAATCTGGCGATCCGAATGCATTATTGAAACACCGAT<br>TTGAGATAATTGAAGGGAGAGACCGAACGATGGCTTGGACAGTGG<br>TAAATAGTATCTGCAACACCACAGGAGTCGATAAGCCTAAATTCC<br>TCCCAGATTTGTATGATTACAAGGAGAACCGATTCATTGAAATTGG<br>AGTGACAAGGAGGGAAGTTCACACATACTACCTAGAAAAGGCAA<br>ATAAGATAAAATCAGAGAAGACACACATTCACATATTCTCATTCA<br>CTGGGGAGGAGATGGCCACCAAAGCTGATTATATCCTTGATGAAG<br>AGAGCAGAGCAAGGATCAAAACCAGGTTGTTCACTATCAGGCAA<br>GAAATGGCCAATAGGGGTCTGTGGGATTCCTTTCGTCAATCTGAGA<br>GAGGCGAAGAGACAATTGAAGAAAGGTTTGAAATCACAGGAACC<br>ATGCGCAGGCTTGCCGACCAAAGTCTCCCACCGAATTTCTCCAGCC<br>TTGAAAATTTTAGAGCCTATGTGGATGGATTCAAACCGAACGGC<br>TGCCTTGAGGGCAAGCTTTCTCAAATGTCAAAAGAAGTGAACGCC<br>AGAATTGAGCCATTCATGAAGACAACACCACGCCCTCTCAGATTA<br>CCTGATGGTCCTCCTTGCTCTCAGCGGTCGAATTCTTACTGATGG<br>ATGCCCTTAAATTGAGCATCGAAGACCCAAGCCATGAGGGAGAAG<br>GTATACCGCTATATGATGCAATCAAATGCATGAAGACGTTTTTTGG<br>TTGGAAAGAGCCCAACATTGTAAAACCACATGTAAAAGGCATAAA<br>TCCCAACTATCTCTTGGCTTGGAAGCAGGTGCTGGTAGAACTCCAA<br>GACATTGAAAATGAAGAGAAATCCCAAAAACAAAAAACATGAA<br>GAAAACAAGCCAACTAAAGTGGGCACTCGGTGAGAATATGGCACC<br>TGAAAAAGTGGACTTTGAGGACTGCAAAGATGTTAGCGATCTAAG<br>ACAGTATGACAGTGATGAACCAGAGCCCAGATCATTATCAAGCTG<br>GATCCAGAGCGAATTCAACAAAGCATGCGAATTGACAGATTCGAG<br>TTGGATTGAACTTGATGAAATAGGAGAAGATGTTGCTCCAATTGA |

| | | | |
|---|---|---|---|
| | | | GCACATTGCGAGTATGAGAAGAAACTACTTCACAGCGGAAGTGTC<br>TCATTGCAGGGCTACTGAATATATAATGAAAGGAGTTTATATAAAT<br>ACAGCCCTGTTGAATTCATCCTGTGCAGCCATGGATGACTTCCAAT<br>TGATTCCAATGATAAGCAAGTGCAGAACCAAAGAAGGAAGACGG<br>AAGACAAATCTATATGGGTTCATTATAAAAGGAAGATCCCATTTG<br>AGGAATGATACCGATGTGGTAAATTTTGTGAGCATGGAGTTCTCTC<br>TTACTGACCCGAGGCTGGAACCACACAAGTGGGAAAAGTACTGTG<br>TTCTCGAAATAGGAGACATGCTCCTACGAACTGCAATAGGCCAAG<br>TATCAAGACCCATGTTTCTTTATGTAAGGACCAATGGGACTTCCAA<br>GATCAAGATGAAATGGGGCATGGAGATGAGGCGATGCCTTCTTCA<br>ATCCCTCCAACAAATTGAGAGCATGATTGAGGCAGAGTCTTCTGTC<br>AAAGAGAAGGACATGACCAAGGAATTCTTTGAAAATAAATCAGAA<br>ACGTGGCCAATTGGGGAATCACCTAAGGGGGTGGAGGAAAGCTCT<br>ATTGGGAAAGTGTGTAGAACATTACTAGCAAAATCTGTATTCAAC<br>AGCCTATATGCATCTCCACAACTTGAGGGGTTTTCAGCTGAGTCGA<br>GAAAGTTACTTCTCATTGTTCAGGCATTTAGGGACAACCTGGAACC<br>TGGGACCTTCGATCTTGGGGGGCTATATGAAGCAATTGAGGAGTG<br>CCTGATTAATGATCCCTGGGTTTTGCTTAATGCATCTTGGTTCAACT<br>CCTTCCTTACACATGCACTGAAATAG |
| UKDD | HA | 1-1704 | **ATG**GAGAAAATAGTGCTTCTTCTTGCAGTGGTTAGCCTTGTTAAAA<br>GTGATCAGATTTGCATTGGTTACCATGCAAACAACTCAACAAAAC<br>AGGTTGACACAATAATGGAAAAAAACGTCACTGTTACACATGCCC<br>AAGACATACTGGAAAAGACACACAACGGGAAGCTCTGCGATCTTA<br>ATGGAGTGAAGCCCCTGATTCTAAAGGATTGTAGCGTAGCTGGGT<br>GGCTCCTTGGAAATCCAATGTGCGACGAGTTCATCAGGGTGCCGG<br>AATGGTCTTACATCGTGGAGAGGGCTAACCCAGCCAACGACCTCT<br>GTTACCCAGGGACCCTCAATGACTATGAGGAACTGAAACACCTAC<br>TGAGCAGAATAAATCATTTTGAGAAAACTCTGATCATCCCCAAGA<br>GTTCTTGGCCCAATCATGAAACATCATTAGGGGTGAGCGCAGCAT<br>GTCCATACCAGGGAGCATCCTCATTTTTTCAGAAATGTGGTATGGCT<br>CATCAAAAAGAACGATGCATACCCGACAATAAAGATAAGCTACAA<br>TAATACCAATCGGGAAGATCTTTTGATACTGTGGGGGATTCATCAT<br>CCCAACAATGCAGAAGAGCAGACAAATCTCTATAAAAACCCAGAC<br>ACTTATGTTTCCGTTGGGACATCAACATTAAACCAGAGATTGGTGC<br>CAAAAATAGCTACTAGATCCCAAGTAAACGGGCAACGTGGAAGAA<br>TGGATTTCTTCTGGACAATTTTAAAACCGAATGATGCAATCCACTT<br>TGAGAGTAATGGAAATTTCATTGCTCCAGAATATGCCTACAAAATT<br>GTCAAGAAGGGGACTCAACAATTATGAAAAGTGAAGTGGAGTAT<br>GGCTACTGCAACACCAAATGTCAAACCCCAATAGGGGCGATAAAC<br>TCTAGCATGCCATTCCACAATATACACCCTCTCACCATCGGGGAAT<br>GCCCCAAATACGTGAAGTCAAACAAATTAGTCCTTGCGACTGGGC<br>TCAGAAATAGTCCTCAAGGGAAAGAAGAAGAAAAAGAGGACTA<br>TTTGGAGCTATAGCAGGGTTTATAGAGGGAGGATGGCAGGGAATG<br>GTAGACGGTTGGTATGGGTACCACCATAGCAATGAGCAGGGGAGT<br>GGGTACGCTGCAGACAAAGAATCCACCCAAAAGGCAGTAGATGG<br>AGTTACCAATAAGGTCAACTCAATCATTGACAAAATGAACACTCA<br>ATTTGAGGCCGTTGGAAGGGAATTTAATAACTTAGAAAGGAGAAT<br>AGAGAATTTAAACAAGAAATGGAAGACGGATTCCTAGATGTCTG<br>GACTTATAATGCTGAACTTTTAGTTCTCATGGAAAATGAGAGAACT<br>CTAGATTTCCATGACTCAAATGTCAAGAACCTTTACGACAAAGTCC<br>GACTACAGCTTAGGGATAATGCAAAGGAGCTGGGTAATGGTTGTT<br>TCGAGTTCTATCACAAATGTGATAACGAATGTATGGAAAGCGTAA<br>GAAATGGGACGTATGACTACCCTAAGTATTCAGAAGAAGCAAGAT<br>TAAAAGAGAAGAAATAAGCGGAGTGAAATTAGAATCAATAGGA<br>ACTTACCAAATACTGTCAATTTATTCAACAGTGGCGAGTTCCCTAG<br>CACTGGCAATCATAGTGGCTGGTCTATCTTTATGGATGTGCTCTAA |

| | | | |
|---|---|---|---|
| | | | TGGGTCGCTACAATGCAGAATTTGCATCTAA |
| UKDD | NP | 1-1497 | **ATG**GCGTCTCAAGGCACCAAACGATCTTATGAACAGATGGAAACTGGTGGAGAACGCCAGAATGCCACTGAAATCAGAGCATCTGTTGGAAGAATGGTTGGTGGAATTGGAAGGTTTTATATACAGATGTGCACTGAACTCAAACTCAGCAATTATGAGGGGAGACTGATCCAGAACAGCATAACAATAGAAAGAATGGTTCTCTCTGCATTTGATGAAAGGAGGAACAAGTACCTGGAAGAACATCCCAGTGCGGGGAAGGACCCAAAGAAAACTGGAGGTCCAATCTACAGAAGAAGAGACGGAAAGTGGATGAGGGAGCTGATTCTGTATGACAAAGAAGAGATCAGAAGGATCTGGCGTCAAGCAAATAATGGAGAAGATGCAACTGCTGGTCTCACCCATCTGATGATCTGGCACTCCAACCTGAATGATGCCACATATCAGAGGACAAGGGCTCTCGTGCGCACTGGAATGGATCCCAGAATGTGCTCTCTGATGCAAGGATCAACTCTCCCAAGAAGGTCTGGAGCTGCTGGTGCAGCAGTAAAAGGGGTCGGAACAATGGTAATGGAATTGATTCGAATGATAAAGCGAGGGATTAATGATCGGAATTTCTGGAGAGGCGAAAATGGAAGAAGGACAAGGATTGCCTATGAGAGAATGTGCAACATCCTCAAAGGGAAATTTCAAACAGCAGCACAAAGAGCAATGATGGATCAAGTGCGAGAAAGCAGGAATCCTGGGAATGCTGAAATTGAAGATCTCATTTTTCTGGCACGTTCTGCACTCATCCTGAGAGGATCAGTGGCCCACAAGTCTTGTCTGCCTGCTTGTGTTTACGGACTTGCTGTGGCCAGTGGATATGACTTTGAGAGAGAAGGATACTCTCTGGTTGGAATAGACCCTTTCCGTCTGCTTCAAAACAGCCAAGTCTTCAGTCTCATTAGACCAAATGAAAACCCAGCACATAAAAGTCAGTTGGTATGGATGGCATGCCATTCAGCAGCGTTTGAGGACCTGAGGGTATCAAGTTTCATCAGAGGGACAAGAGTGGTCCCAAGAGGACAACTATCCACCAGAGGAGTTCAAATTGCATCAAATGAAAACATGGAAACAATGGACTCCAGCACTCTTGAATTGAGAAGCAGATACTGGGCTATAAGAACCAGGAGTGGAGGAAACACCAACCAACAGAGAGCTTCTGCAGGACAAATCAGCGTACAACCCACCTTCTCAGTACAGAGAAATCTTCCCTTTGAAAGAGCGACCATCATGGCGGCATTTACAGGGAACACTGAAGGCAGGACCTCTGACATGAGGACTGAGATCATAAGAATGATGGAAAGTGCCAAACCAGAAGATGTGTCTTTCCAGGGGCGGGGAGTCTTCGAGCTCTCGGACGAAAAGGCAACGAACCCGATCGTGCCTTCCTTTGACATGAGCAACGAAGGATCTTATTTCTTCGGAGACAGTGCAGAGGAGTATGACAATTAA |
| UKDD | NA | 4-1420 | AATCCAAATCAGAAAATAGTAACCATTGGCTCCATTTCATTAGGGTTGGTTGTATTCAATGTTCTACTGCATGCTGTGAGCATCATATTAACAGTGTTAGCCCTGGGGAAGAGTGAAAACAATGGAATCTGCAATGGAACTGTAGTGAGGGAATACAATGAAACAGTTAGAATAGAGAAAGTGACTCAATGGTACAATACTAGCGTAGTCGAATATGTACCGCATTGGAATGAGGGCACTTATATAAATAACACCGAACCAATATGTGATGTCAAGGGCTTTGCACCTTTTTCCAAGGACAATGGGATAAGAGTTGGCTCCAGGGGACATATTTTTGTCATAAGAGAGCCTTTCGTCTCTTGTTCACCTGTAGAGTGCAGGACTTTCTTCCTCACTCAGGGATCTCTACTCAATGACAAACACTCAAATGGAACAGTGAAGGATAGAAGCCCATTCAGAACTCTCATGAGTGTCGAAGTGGGCCAACCACCCAGTGTATATCAAGCCAGGTTTGAAGCTGTGGCATGGTCAGCAACAGCCTGTCATGATGGTAATAAGTGGATGACGATTGGTGTAACAGGGCCAGATTCTAAAGCAGTAGCAGTAGTTCATTACGGAGGGGTGCCTACTGACGTTGTTAACTCCTGGGCAGGAGATATATTAAGAACTCAGGAGTCATCTTGTACTTGCATTCAAGGTAATTGTTATTGGGTAATGACTGACGGTCCTGCCAATAGACAGGCGCAGTATAGAATATACAAAGCAAATCAAGGCAAAATAATTGGCCGAACAGATGTTAGCTTTAGTGGAGGACATATTGAGGAATGTTCTTGTTATCCAAATGATGGTAAAGTGGAATGCGTGTGTAGAGACAACTGGACGGGAACTAACAGGCCTGTACTAATTATTCGCCTGATCTCTCTTACAGGGTTGGGTATTTATGTGCAGGGTT |

| | | | |
|---|---|---|---|
| | | | GCCCAGTGACACTCCAAGAGGGGAAGATACTCAATTTGTCGGTTC<br>ATGCACTAGTCCCATGGGAAATCAGGGATATGGCGTAAAAGGGTT<br>CGGGTTTCGACAGGGAACTGATGTGTGGGTGGGGCGGACAATTAG<br>TCGAACCTCCAGATCAGGATTTGAAATAATAAGGATAAAGAATGG<br>TTGGACGCAAACAAGCAAAGAACAGATTAGAAGACAGGTGGTTGT<br>TGATAACTCGAATTGGTCGGGATACAGTGGGTCTTTCACTTTACCA<br>GTAGAATTGTCTGGGAGGGAATGTTTGGTTCCCTGTTTTTGGGTCG<br>AAATGATCAGAGGTAGGCCAGAAGAGAGAACAATCTGGACCTCTA<br>GTAGCTCCATTGTAATGTGTGGAGTTGATTATGAAATTGCCGATTG<br>GTCATGGCACGATGGAGCTATTCTTCCCTTTGACATCGATAAGACG<br>TAATTTACGA |
| UKDD | MP | -5-982 | GAAAG**ATG**AGTCTTCTAACCGAGGTCGAAACGTACGTTCTCTCTAT<br>CATCCCGTCAGGCCCCCTCAAAGCCGAGATCGCGCAGAGACTTGA<br>AGATGTCTTTGCAGGGAAAAACACCGATCTCGAGGCTCTCATGGA<br>GTGGCTAAAGACAAGACCAATCCTGTCACCTCTGACTAAAGGGAT<br>TTTGGGATTTGTGTTCACGCTCACCGTGCCCAGTGAGCGAGGACTG<br>CAGCGTAGACGCTTCGTCCAGAATGCCCTAAATGGGAACGGGGAT<br>CCAAATAATATGGATAAGGCAGTTAAGCTATATAAGAAGCTGAAA<br>AGAGAGATAACATTCCATGGGGCTAAGGAGGTCGCACTTAGCTAC<br>TCAACCGGTGCACTTGCCAGCTGCATGGGTCTCATATACAACAGG<br>ATGGGAACGGTGACTACAGAAGTGGCTTTTGGCCTAGTGTGTGCC<br>ACTTGTGAGCAGATTGCAGATTCACAGCATCGGTCCCACAGACAG<br>ATGGCAACCATCACCAACCCATTAATCAGACATGAGAACAGAATG<br>GTGCTGGCCAGCACTACAGCTAAGGCCATGGAGCAGATGGCAGGA<br>TCAAGCGAGCAGGCATCAGAAGCCATGGAGGTTGCTAATCAGGCC<br>AGGCAGATGGTACAGGCAATGAGGACAATTGGGACTCATCCTAAC<br>TCTAGTGCTGGTCTGAGAGATAATCTTCTTGAAAATTTGCAGGCCT<br>ACCAGAACCGAATGGGAGTGCAGATGCAGCGATTCAAGTGATCCT<br>CTTGTTGTTGCCGCAAATATCATTGGGATCCTGCACTTGATATTGT<br>GGATCCTTGATCGTCTTTTCTTCAAATGCATTTATCGTCGCCTTAAA<br>TACGGTTTGAAAATAGGGCCTTCTACGGAAGGGGTACCTGAGTCT<br>ATGAGGGAAGAGTACCGGCAGGAACAGCAGAGTGCTGTGGATGTT<br>GACGATGGTCATTTGTCAACATAGAATTGGAGTAA |
| UKDD | NS | -5-849 | ACATA**ATG**GACTCCAACACTGTGTCAAGCTTTCAGGTAGACTGCTT<br>TCTTTGGCATGTCCGCAAACGATTTGCAGACCAAGAACTGGGTGAT<br>GCCCCATTCCTTGACCGGCTTCGCCGAGACCAGAAGTCCCTAAGA<br>GGAAGAGGCAGCACTCTTGGTCTGGACATCGAGACAGCTACTCGT<br>GCGGGAAAGCAAATATTGGAGCGGATTCTGGGGGAAGAATCTGAT<br>GAAACACTTAAAATGAATATTGCTTCTGTACCGACTTCACGCTACC<br>TAACTGACATGACTCTTGAAGAAATGTCAAGAGACTGGTTCATGCT<br>CATGCCCAAGCAGAAAGTAGCAGGTTCTCTCTGCATCAAAATGGA<br>CCAGGCAATAATGGATAAAACCATCATACTGAAAGCAAACTTCAG<br>TGTGATTTTTGATCGGCTGGAAACCCTAATATTACTTAGAGCTTTC<br>ACAGAAGAAGGAGCAATTGTGGGAGAAATCTCACCATTACCTTCT<br>CTTCCAGGACATACTGATGAGGATGTCAAAATTGCAATTGGGGTC<br>CTCATCGGAGGGCTTGAATGGAATGATAACACAGTTCGAGTCTCT<br>GAAACTCTACAGAGATTCACTTGGAGAAGCAGTAATGAGGATGGG<br>AGACCTTCACTCCCTTCAAAACAGAAACGGAAAATGGCGAGAACA<br>ATTGAGTCAGAAGTTCGAGGAAATAAGATGGCTGATTGAGGAAAT<br>GCGACATAGATTGAAGATCACAGAGAACAGCTTCGAACAAATAAC<br>GTTTATGCAAGCTTTACAACTATTGCTTGAAGTGGAGCAAGAGATA<br>AGAACCTTCTCGTTTCAGCTTATTTAATGATAA |

| Sample | SP | Platform | Method | Reads | Nucleotides | Influenza reads | Influenza Nucleotides |
|---|---|---|---|---|---|---|---|
| **DETU** | 1 | Illumina MiSeq | RNA-Seq+PCR | 35,397,942 | 4,768,436,983 | ca. 21,238,765 | ca 2,861,062,190 |
| | 3 | 454 | Amplicon | 78,028 | 25,829,288 | 75,913 | 25,692,541 |
| | 2 | Illumina MiSeq | RNA Shot gun | 1,394,424 | 417,805,080 | 1,062,401 | 318,461,282 |
| **NLCH** | 1 | Illumina MiSeq | RNA-Seq+PCR | 45,091,902 | 6,487,449,580 | 1,454,528 | 203,647,299 |
| | 3 | 454 | Amplicon | 32,661 | 12,458,090 | 32,661 | 12,458,090 |
| | 2 | Illumina MiSeq | RNA Shot gun | 1,148,978 | 344,137,436 | 373,742 | 112,011,370 |
| **UKDD** | 1 | Illumina MiSeq | RNA-Seq+PCR | 10,214,524 | 768,562,277 | 867,355 | 64,794,700 |
| | 3 | 454 | Amplicon | 49,993 | 18,897,160 | 48,769 | 18,821,757 |
| | 2 | Illumina MiSeq | RNA Shot gun | 1,512,512 | 421,870,650 | 1,039,962 | 294,863,446 |

**Software**: CLC Genomics Workbench 8

Black = applied for SP1,SP2 and SP3 data
<span style="color:blue">Blue = applied for SP3 data</span>
<span style="color:green">Green = applied for SP1 and SP2 data</span>


1. **Demultiplex**

   File: JGJ0HAZ01.sff
   Define tags:
        Barcode length: 11
        Sequence: 1-1000
   Set barcode options:
        Search both strands = yes
        Barcodes: MID sequences Roche (1-6)
   Result handling:
        Create list of reads without barcode = yes
        Create report = no
        Save = yes


1. **Workflow Map reads to reference beta**
   a. **Trim sequences**
        Quality trimming:
             Ambiguous trim = Yes
             Ambiguous limit = 2
             Quality trim = Yes
             Quality limit = 0,05 = phred score = 20
        Adapter trimming:
             <span style="color:blue">Trim adapter list = NA</span> <span style="color:green">Trim adapter library 2</span>
             Use colorspace = No
             Search on both strands: Yes
        Sequence filtering:
             Remove 5' terminal nucleotides = <span style="color:blue">No</span> <span style="color:green">Yes</span>
             Number of 5' terminal nucleotides = <span style="color:blue">NA</span> <span style="color:green">30</span>
             Remove 3' terminal nucleotides = <span style="color:blue">No</span> <span style="color:green">Yes</span>
             Number of 3' terminal nucleotides = <span style="color:blue">NA</span> <span style="color:green">30</span>
             Discard short reads = Yes
             Minimum number of nucleotides in reads = 15
             Discard long reads = Yes
             Maximum number of nucleotides in reads = 1.000
        Result handling:
             Save discarded sequences = Yes
             Save broken pairs = No
             Create report = Yes
             Result handling: Save

➔ SP3 data files saved as *.fastq and shared via DataHub


   b. **Map reads to reference data**

Select sequencing reads
  Trimmed reads. For 454 data enter both MID's.
References
  Use reference files mentioned above. Consensus sequences per sample derived
  from the consensus sequences of the different institutions
  Masking mode = No masking
  Exclude annotated = NA
  Include annotated only = NA
Mapping options:
  Mismatch cost = 2
  Cost of insertions and deletions = Affine gap cost
  Insertion open cost = 7
  Insertion extend cost = 2
  Deletion open cost = 7
  Deletion extend cost = 2
  Length fraction = 0,7
  Similarity fraction = 0,9
  Global alignment = No
  Non-specific match handling = Map randomly
Results handling:
  Output mode = Create stand-alone read mappings
  Create report = Yes
  Collect un-mapped reads = Yes
  Save = yes


2. **Workflow Realign and detect variants**
  Local realignment
      Realignment settings:
          Realign unaligned ends = Yes
          Multi-pass realignment = 2
          Guidance-variant track = Not set
      Result handling:
          Output mode = Create reads track
          Output track of realigned regions = No


  Indels and Structural variants
      Select read mappings
          Locally realigned file
      Select settings
          P-Value threshold = 0,0001
          Maximum number of mismatches = 3
          Filter variants = Yes
          Minimum number of reads = 2
          Reference masking: NA
      Result handling:
          Create report = No
          Create breakpoints = No
          Create InDel variants = Yes
          Create structural variations = No
          Save = yes


  Local realignment

Realignment settings:
    Realign unaligned ends = Yes
    Multi-pass realignment = 2
    Guidance-variant track = Locally realigned (InDel)-file
    Force realignment to guidance-variants = No
Result handling:
    Output mode = Create reads track
    Output track of realigned regions = No

Low frequency variant detection
    Select read mappings
        Locally realigned – locally realigned files
    Low frequency variant parameters
        Required significance (%) = 1,0
    General filters
        Ignore positions with coverage above = 100.000
        Restrict calling to target regions = Not set
        Ignore broken pairs = Yes
        Ignore non-specific matches = Reads
        Minimum coverage = 2
        Minimum count = 2
        Minimum frequency (%) = 1,0
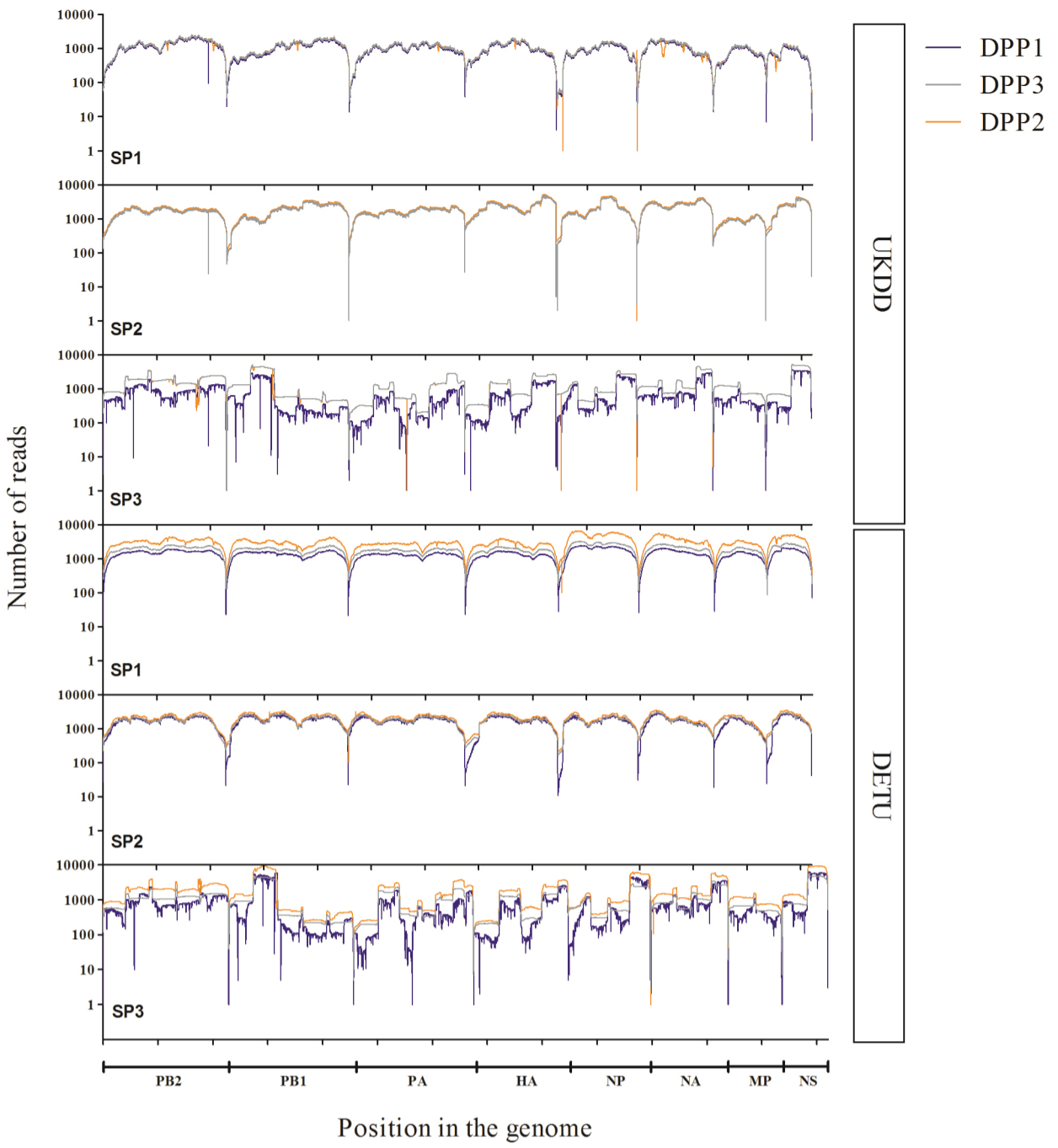    Noise filters
        Base quality filter = Yes
        Neighborhood radius = 5
        Minimum central quality = 0
        Minimum neighborhood quality = 0
        Read direction filter = Yes
        Direction frequency (%) = 5,0
        Relative read direction filter = Yes
        Significance (%) = 1,0
        Read position filter = Yes
        Significance (%) = 1,0
        Remove pyro-error variants = No (454 data checked with and
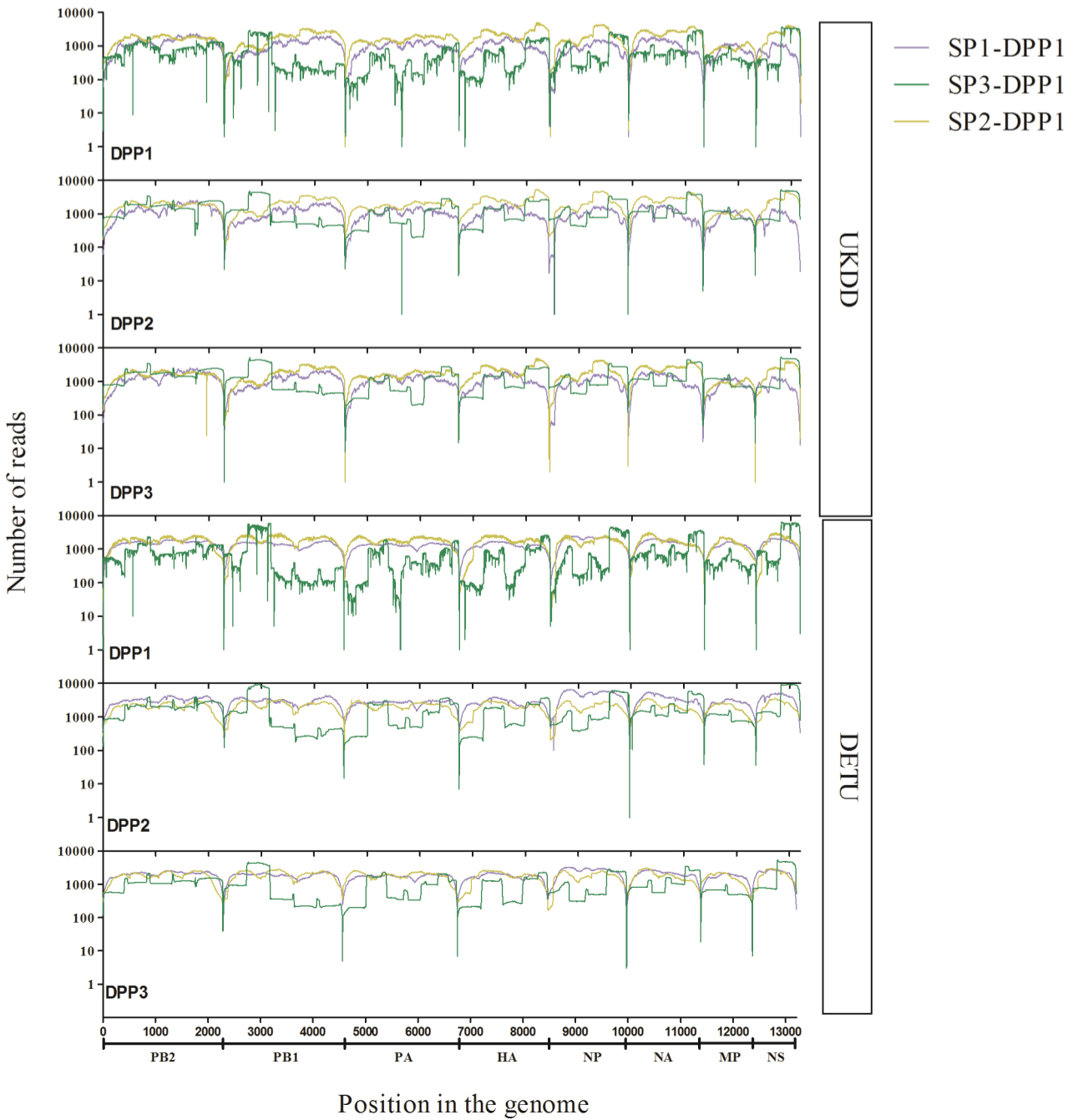            without, no difference for mSNV identification)
    Result handling:
    Create track = Yes
    Create annotated table = Yes
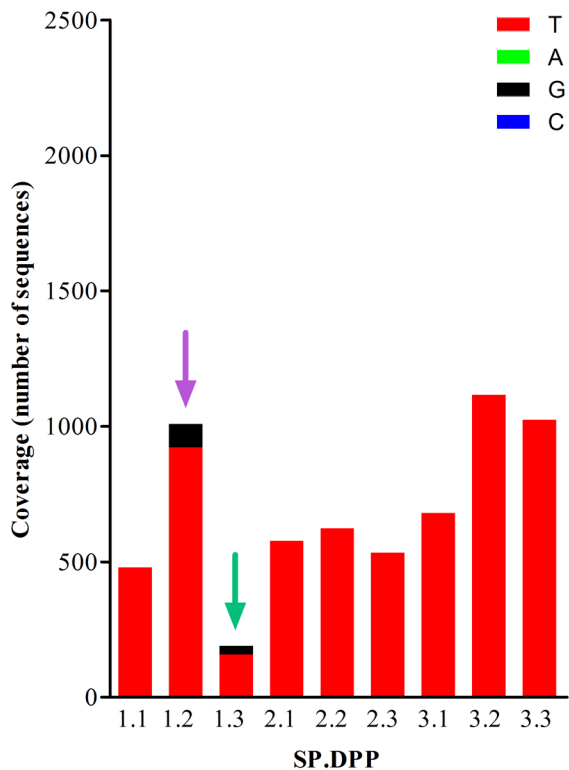    Create report = No

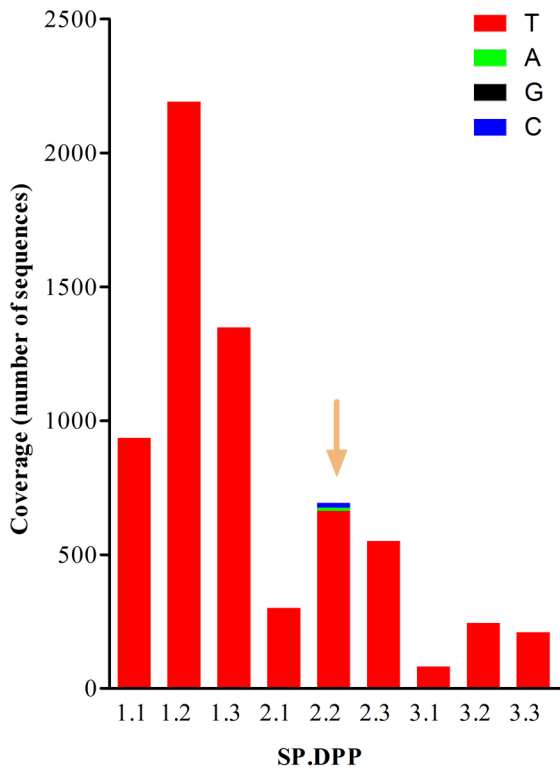**A: DPP derived differences: DPP analyses results per SP**

DPP1
DPP3
DPP2

UKDD

DETU

SP1
SP2
SP3
SP1
SP2
SP3

Number of reads

PB2  PB1  PA  HA  NP  NA  MP  NS

Position in the genome

**B: SP derived differences: SP datasets analysed per DPP**

SP1-DPP1
SP3-DPP1
SP2-DPP1

UKDD

DETU

DPP1
DPP2
DPP3
DPP1
DPP2
DPP3

Number of reads

0  1000  2000  3000  4000  5000  6000  7000  8000  9000  10000  11000  12000  13000

PB2  PB1  PA  HA  NP  NA  MP  NS

Position in the genome