

1 **Reliable variant calling during runtime of Illumina sequencing**

2 Tobias P. Loka¹, Simon H. Tausch^{1,2,3}, Bernhard Y. Renard^{1*}

3 ¹ Bioinformatics Division (MF 1), Department for Methods Development and Research Infrastructure

4 ² Centre for Biological Threats and Special Pathogens: Highly Pathogenic Viruses (ZBS 1)

5 ³ Present address: German Federal Institute for Risk Assessment (BfR), Department of Biological Safety, Berlin,

6 Germany

7 * To whom correspondence should be addressed

8 **Abstract**

9 The sequential paradigm of data acquisition and analysis in next-generation sequencing
10 leads to high turnaround times for the generation of interpretable results. We designed a
11 workflow using an advanced real-time read mapping approach to obtain reliable variant calls
12 for human whole-exome data still during the sequencing process. When compared to
13 standard routines, our live variant calling approach enables considerably faster interventions
14 in clinical applications such as pathogen characterization and the determination of drug
15 resistances in disease outbreaks or the design of individually tailored vaccines in precision
16 medicine. Besides variant calling, our approach can be adapted for a plethora of other
17 mapping-based analyses.

18 **Keywords**

19 Next Generation Sequencing, variant calling, real-time analysis, short read alignment,
20 Illumina sequencing, pathogen identification, personalized medicine, drug resistance

21 **Background**

22 Common workflows for the analysis of next-generation sequencing (NGS) data can only be
23 applied after sequencing has finished. In time-critical applications, however, this sequential
24 paradigm of data acquisition and analysis is one of the main bottlenecks leading to high
25 turnaround times. Examples for such time-critical analyses range from the production of

26 individually tailored vaccines for cancer immunotherapy [2], to the determination of *M.*
27 *tuberculosis* drug resistances [3], and to the identification of pathogens, virulence factors,
28 drug resistances and paths of disease transmission in infectious disease outbreaks [4, 5].
29 While having considerably higher turnaround times than alternative approaches such as
30 molecular tests, NGS provides a more open view as well as more extensive and reliable
31 results. During bioinformatics analysis of NGS data, variant calling is a crucial step to find
32 differences in the genomic sequence of the investigated sample when compared to a
33 reference genome. Thereby, valuable information for the treatment of a patient can be
34 obtained such as strain level classification and drug resistances of a pathogen or individual
35 characteristics of healthy and defected tissue. To reduce the turnaround time for the
36 generation of NGS-based results and to enable fast and accurate treatment of patients, we
37 designed a workflow to obtain variant calling results before sequencing has finished. The
38 workflow is based on real-time read mapping results with HiLive2 followed by fast and
39 accurate variant calling with xAtlas. In doing so, live results can be obtained several hours
40 before all data are written by the sequencer and provide increasing insights into the sample
41 over sequencing time.

42 **Methods**

43 In this chapter, we provide a description of our workflow and the methods for evaluation.
44 More detailed information for each step such as direct links to the data and the executed
45 commands can be found in the supplementary material. All software versions used in this
46 study are listed in the *Software versions* section (Table 1).

47 **Implementation of HiLive2**

48 HiLive2 is based on a novel algorithm using the efficient FM-index implementation of the
49 SeqAn library [6]. Each alignment starts with a short exact match sequence of length k ,
50 referred to as seed k -mer. The length of this k -mer is set according to the size of the
51 reference genome and the read length or can be manually specified by the user. Starting

52 from the seed k -mer, the alignment is extended in the sequencing direction of the read. Apart
53 from a front soft clip that occurs due to k -mer mismatches at the beginning of a read, the
54 algorithm guarantees to produce optimal mapping results for non-affine gap costs. HiLive2 is
55 an all-mapper by design meaning that all alignments up to a specified score can be found.
56 However, the default output option of HiLive2 is to return only one best alignment for each
57 read which is the expected behavior for most analyses. During output, temporary files are
58 stored for all output cycles such that new output with the same or different output options can
59 be created by a separate executable. Output is created in the well-established BAM or SAM
60 format.

61 **Data download and conversion**

62 The human reference genome hg19 was obtained from NCBI and only considered
63 chromosomes 1-22, X and Y. Alternative regions were omitted. The sequences were stored
64 in a single multi-FASTA file. For the evaluation of variant calls with RTG Tools [7], the
65 reference genome was converted to SDF format.

66 Seven whole-exome sequencing (WES) data sets of the human individual NA12878 were
67 downloaded from EBI in FASTQ format. For read mapping with HiLive2, read pairs were
68 converted to Illumina base call file format (BCL), distributed on one lane and 64 tiles. There
69 were four different definitions for exome capture region definition required for the different
70 data sets (cf. Table 2). The regions were obtained in BED format from the respective
71 producer, if available. Whenever multiple definition files were provided, the primary target
72 regions were selected.

73 Gold standard variants for the individual NA12878 were downloaded from the Genome in a
74 Bottle (GIAB) consortium [8] and regularized with the `vcfallelicprimitives` tool of VCFtools [9].
75 SNPs and indels of the gold standard were stored in two separated files and filtered out
76 against the exome capture regions using `Bedtools` [10] `intersect`. The resulting files were
77 used as the gold standard for data sets using the respective exome capture definition. During

78 the evaluation of the results, only variant calls in high confidence homozygous regions which
79 were obtained from GIAB were considered.

80 **Real-time read mapping with HiLive2**

81 The index of human reference genome hg19 for HiLive2 was built with default parameters.
82 The creation of base call files by the sequencing machine was simulated using a script for
83 sequencing simulation with a sequencing profile for HiSeq2500 machines in rapid mode and
84 using dual barcodes. As no barcodes were present in our data sets, no data was written by
85 the sequencing simulator for the respective cycles. HiLive2 was run in fast mode allowing
86 faster turnaround times at the expense of slightly lower recall. Technical parameters as
87 lanes, tiles and read length were set for each run according to the respective data sets. In
88 general, we chose cycles 30, 40, 55, 75 and 100 for each of the two reads as output cycles.
89 For data sets with read lengths other than 2 x 100bp, we adapted the output cycle numbers
90 to 30, 40 and 50 (SRR292250) or 30, 40, 55 and 76 (SRR098401). We used the
91 recommended number of threads (1 thread per tile) for HiLive2 resulting in 64 threads for all
92 data sets.

93 **Read mapping with Bowtie 2**

94 The index of the human reference genome hg19 for Bowtie 2 was built with default
95 parameters. Read mapping with Bowtie 2 was performed with default parameters using 10
96 threads.

97 **Variant calling with xAtlas**

98 Variant calling with xAtlas was performed for each chromosome individually. Therefore, the
99 alignment files of HiLive2 or Bowtie 2 were split in 24 files (one for each chromosome). The
100 resulting files were sorted and indexed using samtools [11]. Afterwards, variants were called
101 with xAtlas for the respective exome capture regions using default parameters. Sorting,
102 indexing and variant calling was performed with 24 threads (one thread per chromosome).

103 The resulting VCF files were merged using VCFLIB vcf-concat [12] for SNVs and indels
104 separately.

105 **Measure of turnaround time**

106 The sequencing simulation script provides timestamps for each written sequencing cycle.
107 These timestamps were compared to the system time stamps for the last modification of the
108 alignment output files of HiLive2. The time span between both time stamps describes the
109 alignment delay of HiLive2. Additionally, we measured the clock time of the xAtlas pipeline.
110 The sum of sequencing time until the respective cycle, the alignment delay of HiLive2 and
111 the clock time of xAtlas yields the overall turnaround times of our workflow.

112 **Evaluation with RTG Tools**

113 We used the vcfeval program of RTG Tools for the validation of variant calling results. We
114 used the gold standard for the respective data set (depending on the used exome kit) as
115 baseline and the variant calling output of xAtlas or GATK as call. The human reference hg19
116 in SDF format was used as reference template. Only variant calls being included in the high-
117 confidence regions for individual NA12878 provided by the GIAB consortium were
118 considered for validation. We ran RTG Tools with 24 threads and used the squash-ploidy
119 and all-records parameters. For variant calls produced by xAtlas, we additionally defined
120 QUAL as the field for variant call quality. For GATK, the GQ field is chosen by default. RTG
121 Tools vcfeval returns a list of statistical measures for different thresholds of the variant call
122 quality field, including precision and recall. These values served as input for the precision-
123 recall curves shown in Supplementary Fig. 1 and used for the calculation of the area under
124 the precision-recall curves (APR).

125 **Software versions**

126 All analyses performed in this study were done with the software versions listed in Table 1.

127 **Table 1** List of software used in this study. Software with source Bioconda was installed with the environment
128 management software conda (<https://conda.io>) and obtained from the Bioconda channel [13].

Name	Version	Source	Used for
Bedtools [10]	2.21.0	https://github.com/arq5x/bedtools2	Vcf and Bed file processing
Bowtie 2 [14]	2.3.4.1	Bioconda	Read alignment
GATK [15]	3.8	Bioconda	Alignment file processing
HiLive2	2.0	https://gitlab.com/rki_bioinformatics/hilive2	Real-time read alignment
RTG Tools [7]	3.9	Bioconda	Benchmark of variant calls
Samtools [11]	1.8	Bioconda	SAM/BAM file processing
VCFLIB [12]	1.0.0_rc1	Bioconda	Vcf file processing
VCFtools [9]	0.1.12	https://sourceforge.net/projects/vcftools/	Vcf file processing
xAtlas [16]	0.1	Bioconda	Fast variant calling

129 Results

130 Implementation and experimental setup

131 To allow for faster NGS-based diagnosis and treatment, we developed a workflow to produce
132 high-quality variant calls based on intermediate read mapping results while sequencing is still
133 running. This approach allows reliable and fast variant calling results without reducing the
134 final sequencing coverage or quality. Therefore, we adapted the real-time read mapper
135 HiLive [17] that gave output at the end of sequencing, using a novel algorithm based on an
136 efficient FM-index implementation for continuously analyzing sequencing results during
137 runtime. The new version (HiLive2) achieves scalability to larger indices such as the
138 complete human reference genome. At the same time, the algorithm comes with improved
139 performance in terms of runtime, memory and data storage and overcomes heuristic
140 elements that were present in previous version of HiLive. The high scalability and accuracy
141 enable the combination of real-time read mapping results with complex follow-up analyses.
142 To demonstrate the power of such analyses, we performed variant calling on seven WES
143 data sets of the human individual NA12878 from the CEPH Utah Reference Collection (cf.
144 Table 2) using the real-time read mapping results of HiLive2 as input data.

145 **Table 2** Summary of data sets evaluated in this study. Information about sequencing platform, exome capture and
146 coverage were adopted from Hwang et al. (2015) [1].

Accession No.	Platform	Exome capture	Exome coverage	Reads ¹	Read length
SRR098401	HiSeq2000	SureSelect v2	116.84x	114M	2 x 76bp
SRR292250	HiSeq2000	SeqCap EZ v2	116.06x	85M	2 x 50bp
SRR515199	HiSeq2000	SureSelect v4	298.45x	167M	2 x 100bp
SRR1611178	HiSeq2000	SeqCap EZ v3	79.93x	45M	2 x 100bp
SRR1611179	HiSeq2000	SeqCap EZ v3	79.84x	45M	2 x 100bp
SRR1611183	HiSeq2500	SeqCap EZ v3	129.94x	74M	2 x 100bp
SRR1611184	HiSeq2500	SeqCap EZ v3	111.90x	64M	2 x 100bp

147 ¹ M:= millions

148 For variant calling, we used the fast variant caller xAtlas [16] which shows comparable
149 accuracy to established methods at much lower runtime. We compared our results to read
150 mapping with Bowtie 2 [14] and variant calling with either xAtlas or GATK HaplotypeCaller
151 [15] for the same data sets. Accuracy was determined by comparing the results to the well-
152 established high-confident variant calls for the human individual NA12878 published by the
153 Genome in a Bottle (GIAB) consortium [8]. As benchmarking method we used the area under
154 the precision-recall curve (APR).

155 **Accuracy of real-time results**

156 In Illumina sequencing, all reads are sequenced in parallel. In each so-called sequencing
157 cycle, sequence information of one additional nucleotide is obtained for all reads. Thus, the
158 current length of a read equals the number of the respective cycle (e.g., 40 nucleotides after
159 cycle 40). To demonstrate the capability of our approach to provide interpretable results
160 during runtime, we applied our workflow at different stages of sequencing. We expected our
161 live results to show higher accuracy for higher cycles due to the increasing amount of
162 available sequence information. At the same time, we analyzed whether the detected
163 variants in early sequencing cycles are as reliable as variants called at the end of
164 sequencing. This is a crucial criterion for our real-time workflow since interpretation of live
165 results is only meaningful when based on reliable variant calls. Therefore, besides comparing

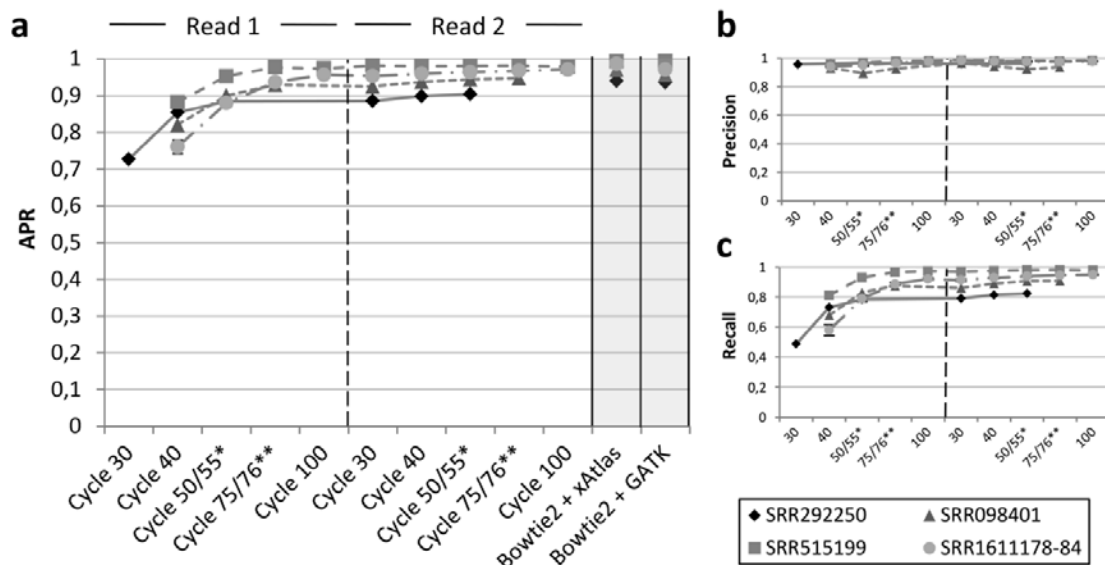


Fig. 1: Area under a precision-recall curve (APR) for SNP calling in seven data sets at different sequencing cycles. SNP calling was performed with xAtlas using real-time read mapping results of HiLive2. Results for the samples SRR1611178, SRR1611179, SRR1611183 and SRR1611184 were combined to a single data series (SRR1611178-84). Error bars for this data series show the standard deviation. The vertical, ticked line in the middle of the plot divides the first and second read. **a** The gray columns show APR values using Bowtie 2 for read mapping and xAtlas (left) and GATK-HC (right) for variant calling. The data for Bowtie 2 + GATK were taken from Hwang et al. (2015) [1]. The real-time workflow with HiLive2 and xAtlas provides first results after 40 sequencing cycles (30 cycles for SRR292250). An APR value greater than 0.9 is reached after 75 cycles for all data sets with a minimal read length of 75bp. Until end of sequencing, there is a moderate increase of the APR. **b** Precision with a quality threshold of 1 for variant calling with xAtlas. The results show no precision lower than 0.89 for all sequencing cycles. This indicates that results in early sequencing cycles are already reliable. **c** Recall with a quality threshold of 1 for variant calling with xAtlas. The results show strong improvements from the first results available until the end of the first read. The progression of all curves is similar to that of the APR curve (cf. Fig. 1a), indicating the correlation between those two measures. *Cycle 50 for SRR292250, cycle 55 for all other data sets. **Cycle 76 for SRR098401, cycle 75 for all other data sets.

166 the APR values of different sequencing cycles, we also examined precision and recall
 167 separately. Fig. 1a shows the progression of the APR values for SNP calling in all analyzed
 168 data sets with increasing sequencing time. In cycle 30, sequence information was not
 169 sufficient to call any variants with the given parameter settings for six of seven data sets. For
 170 data set SRR292250, read mapping parameters were adapted by HiLive2 automatically due
 171 to the short read length of 50bp. This led to earlier results after 30 cycles while first results
 172 were available after cycle 40 for all other data sets. Results show a continuous increase of
 173 the APR values for all cycles of the first read. In cycle 75, an APR larger than 0.9 was
 174 achieved for all data sets with sufficient read length. Afterwards, the APR values continue

175 increasing moderately. When regarding the progression of precision (Fig. 1b) and recall (Fig.
176 1c) over sequencing time separately, it can also be observed that lower APR values for
177 earlier sequencing cycles are mainly caused by a lower recall while precision changes only
178 slightly with more sequence information. The same conclusions are supported by the
179 individual precision-recall curves for all data sets which show a large increase of the recall
180 but only minor changes of specificity over sequencing time (cf. Supplementary Fig. 1). This
181 indicates that live results are highly reliable and can therefore serve for early interpretation
182 and problem-specific follow-up analyses. The increasing number of SNP calls in subsequent
183 cycles provides additional information for complementing the previous interpretation of the
184 data. However, the final results with HiLive2 show slightly lower maximum recall values than
185 the same workflow applied to read mapping results of Bowtie 2 (cf. Supplementary Fig. 1).
186 This can be explained by the read mapping approach of HiLive2 which tolerates only a
187 specified number of errors for a read. Thus, regions with a high number of variations may be
188 lowly covered which leads to undetected variants. The same effect is somewhat stronger for
189 indels as the algorithm only tolerates indels with a maximum length of three nucleotides by
190 default. While this behavior led to a lower recall than based on read mapping with Bowtie 2,
191 the results showed comparable or higher precision (cf. Supplementary Fig. 1). Thus,
192 although focussing on SNPs in this study, our workflow can also provide valuable insights
193 about small indels.

194 **Turnaround time of the workflow**

195 Besides the accuracy of results, turnaround time is the second crucial factor for NGS-based
196 real-time analyses. Thereby, live results should be available as soon as possible after the
197 data of the respective sequencing cycle was written without showing significant delay in any
198 stage of sequencing.

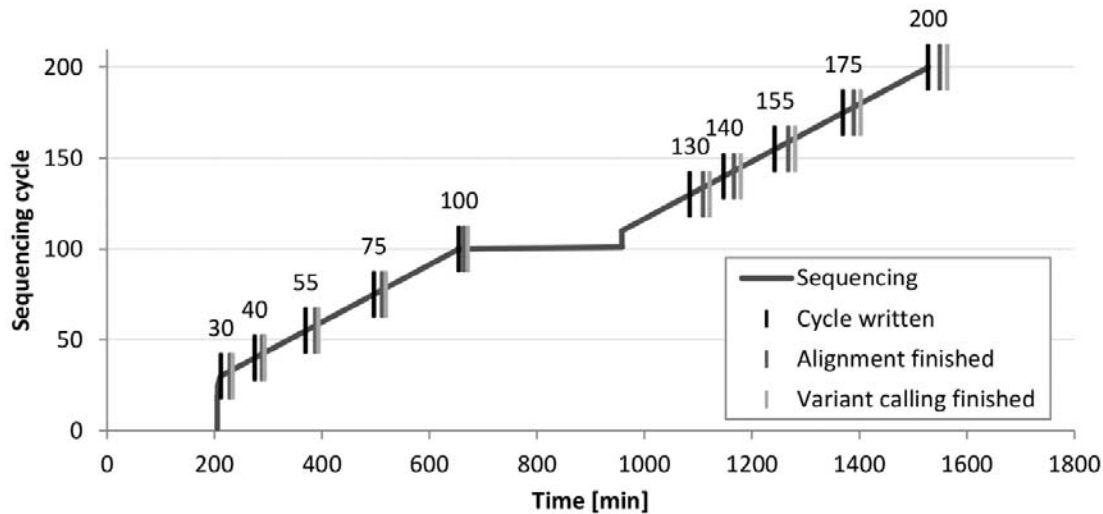


Fig. 2: Turnaround time of our workflow for data set SRR1611178. For each cycle, the first vertical line indicates the time point when the data for the respective cycle was completely written. The second line shows when the alignment output of HiLive2 is written. The third line indicates the end of our workflow resulting in the output of variant calls for the respective cycle. After cycle 100, there is no increase in the sequencing cycle for several hours due to the time spent for sequencing indices and initialization of the second read. In average, the time from data written by the sequencing machine and the output of variant calls is lower than 21 minutes for the first read and around 35 minutes for the second read. Final results are written 36 minutes after sequencing finished.

199 We measured the turnaround time of real-time mapping with HiLive2 and subsequent variant
200 calling with xAtlas for the same runs that delivered the accuracy results shown before. All
201 computations were run on a 128-core machine (Intel® Xeon® CPU E5-4667 v4 @
202 2.20 GHz, 45 M Cache) with 500GB RAM, using a maximum of 65 threads per data set.
203 Fig. 2 shows an overview for the turnaround time of our workflow for different sequencing
204 cycles for data set SRR1611178. On average, variant calling results were available less than
205 half an hour after the data were written by the sequencing machine. The results for cycle 40
206 were written after 294 minutes, showing that first reliable and interpretable variant calls were
207 available less than five hours after sequencing started. For higher coverage data sets, such
208 as data set SRR515199, the alignment delay of HiLive2 and runtime of xAtlas increase when
209 performed with the same number of threads (cf. Supplementary Fig. 2). However, five of the
210 seven data sets in this study showed a maximum time span of less than one hour from data
211 output to interpretable results for each sequencing cycle. The turnaround times for all data
212 sets are shown in Supplementary Fig. 2.

213 **Discussion**

214 Our results show that real-time read mapping results in very early stages of sequencing can
215 already serve as input for variant calling and deliver confident results. However, the quantity
216 of analysis results (i.e. the number of called variants in this study) increases with a growing
217 number of sequenced nucleotides per read. Live analyses can therefore provide first relevant
218 insights into the data while the analysis becomes more comprehensive with ongoing
219 sequencing. Thereby, our approach does not only apply to the presented use case of variant
220 calling. We rather see the enormous potential of real-time read mapping to provide means for
221 a wide range of complex follow-up analyses.

222 In clinical applications and infectious disease outbreaks, the turnaround time of analyses is a
223 critical factor for an effective treatment of patients. However, a high depth of analysis and an
224 open perspective for unexpected findings are further crucial criteria in such scenarios.
225 Despite its significantly higher turnaround times than alternative methods, NGS presents an
226 established analysis method in several time-critical applications due to its high sensitivity. For
227 example, a comprehensive report of vancomycin resistant *Enterococcus faecium* infections
228 in three patients was created in 48.5 hours including over-night culturing using an Illumina
229 MiSeq benchtop sequencer [18]. A different study reports a workflow for *M. tuberculosis*
230 diagnostics including phylogenetic placement being finished in 44h on an Illumina MiSeq or
231 16h on an Illumina MiniSeq [19]. To further accelerate such analyses, we showed that live
232 results can deliver a major proportion of the full analysis depth already in a fraction of the
233 final sequencing time. These results demonstrate the enormous potential of our approach to
234 reduce the turnaround time from sample arrival to meaningful analysis output by several
235 hours. However, alternative approaches can also be highly valuable for different scenarios.
236 Molecular approaches are usually highly reliable and provide answers to specific questions in
237 a very short timeframe and at much lower costs. For example, the detection of 25 genetic
238 mutations in *M. tuberculosis* that confer to drug resistances can be finished in approximately
239 two hours with a variation of the molecular GeneXpert test [3]. Even when providing live

240 results, such short turnaround times are currently not feasible with NGS-based approaches
241 due to the required time for sample preparation. However, NGS enables more detailed and
242 unbiased analyses ranging from strain level identification to the determination of infection
243 chains. Another interesting technology for time-critical applications is nanopore sequencing.
244 It was shown that metagenomic detection of viral pathogens can be achieved in less than six
245 hours [20]. While nanopore sequencing shows a high portability as an additional benefit, this
246 and other current long-read technologies are still expensive and limited by their
247 comparatively low coverage and high error rates. It is therefore hard to reliably identify lowly
248 abundant pathogens, genetic variants, parallel infections or the presence of viral
249 quasispecies. Thus, especially when it comes to these or other questions going beyond the
250 identification of highly abundant pathogens in time-critical applications, real-time analyses for
251 Illumina sequencing can be of great benefit.

252 **Conclusion**

253 We consider our new real-time workflow for NGS to be a complementary method to
254 molecular tests and ultra-portable, long-read sequencing for time-critical analyses. It fills the
255 current gap of short turnaround times, an open-view perspective and high sequencing
256 coverage which is essential for a plethora of applications such as pathogen identification and
257 characterization, personalized vaccine design or epidemiological analyses. Therefore, we are
258 convinced that our approach will improve the ability for fast interventions in exceptional
259 clinical situations, personalized medicine and infectious disease outbreaks.

260 **Abbreviations**

261 **APR:** area under the precision-recall curve

262 **NGS:** next-generation sequencing

263 **WES:** whole exome sequencing

264 **Declarations**

265 **Acknowledgements**

266 We thank Andrea Thürmer and Aleksandar Radonić for their input concerning technical
267 aspects of Illumina sequencing. We further thank Wojciech Dabrowski for infrastructural
268 support and Thilo Muth for critical reading of the manuscript and his highly valuable
269 suggestions. We also thank Martin Lindner, Jakob Schulze, Benjamin Strauch and Kristina
270 Kirsten for their work on HiLive.

271 **Funding**

272 This work was supported by the German Federal Ministry of Health [IIA5-2512-FSB-725 to
273 B.Y.R. and 2515NIK043 to S.H.T.].

274 **Availability of data and materials**

275 The source code of HiLive2 is published under BSD-3-clause license and available for public
276 download on https://gitlab.com/rki_bioinformatics/hilive2. It comes with extensive
277 documentation and sample data. Scripts that were used for our analyses are provided as
278 Supplementary Software.

279 Sequencing data of the individual NA12878 is publicly available on the NCBI Short Read
280 Archive and on the EBI FTP server. Gold standard variant calls are publicly available from
281 the Genome in a Bottle Consortium. Human reference genome hg19 was obtained from the
282 NCBI FTP server. Exome capture targets are available from the manufacturers or from third
283 party resources.

284 **Authors' contributions**

285 B.Y.R. and T.P.L. conceived the study. T.P.L. performed the implementation of HiLive2.
286 S.H.T. continuously supported the development of HiLive2 and evaluated the performance of
287 HiLive2 on different types of data. T.P.L. designed the workflow and performed the analyses.
288 All authors were involved in the preparation of the manuscript and approved the final version.

289 **Ethics approval and consent to participate**

290 Not applicable.

291 **Consent for publication**

292 Not applicable.

293 **Competing interests**

294 The authors declare no competing interests.

295 **References**

296 1. Hwang S, Kim E, Lee I, Marcotte EM. Systematic comparison of variant calling
297 pipelines using gold standard personal exome variants. *Scientific reports*
298 2015;5:17875.

299
300 2. Sahin U, Tureci O. Personalized vaccines for cancer immunotherapy. *Science*
301 2018;359:1355-1360.

302
303 3. Rubin EJ. TB diagnosis from the Dark Ages to fluorescence. *Nature microbiology*
304 2018;3:268-269.

305
306 4. Quainoo S, Coolen JPM, van Hijum S, Huynen MA, Melchers WJG, van Schaik W,
307 Wertheim HFL. Whole-Genome Sequencing of Bacterial Pathogens: the Future of
308 Nosocomial Outbreak Analysis. *Clinical microbiology reviews* 2017;30:1015-1063.

309
310 5. Gilchrist CA, Turner SD, Riley MF, Petri WA, Jr., Hewlett EL. Whole-genome
311 sequencing in outbreak analysis. *Clinical microbiology reviews* 2015;28:541-563.

312
313 6. Doring A, Weese D, Rausch T, Reinert K. SeqAn an efficient, generic C++ library for
314 sequence analysis. *BMC Bioinformatics* 2008;9:11.

315
316 7. Cleary JG, Braithwaite R, Gaastra K, Hilbush BS, Inglis S, Irvine SA, Jackson A, Littin
317 R, Rathod M, Ware D *et al.* Comparing Variant Call Files for Performance
318 Benchmarking of Next-Generation Sequencing Variant Calling Pipelines. 2015;
319 doi:10.1101/023754.

320
321 8. Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, Salit M. Integrating
322 human sequence data sets provides a resource of benchmark SNP and indel
323 genotype calls. *Nature biotechnology* 2014;32:246-251.

324

- 325 9. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,
326 Lunter G, Marth GT, Sherry ST *et al.* The variant call format and VCFtools.
327 Bioinformatics 2011;27:2156-2158.
- 328
329 10. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
330 features. Bioinformatics 2010;26:841-842.
- 331
332 11. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
333 Durbin R. The Sequence Alignment/Map format and SAMtools. Bioinformatics
334 2009;25:2078-2079.
- 335
336 12. Garrison E, Yandell M, Shapiro M, Marth G, Durbin R, Eichler EE, Consortium TV,
337 Kronenberg ZN. VCFLIB: an ensemble of methods for variant manipulation and
338 population genetics. 2016. <https://github.com/vcflib/vcflib>. Accessed 19.05.2018.
- 339
340 13. Dale R, Grüning B, Sjödin A, Rowe J, Chapman BA, Tomkins-Tinch CH, Valieris R,
341 Köster J. Bioconda: A sustainable and comprehensive software distribution for the life
342 sciences. bioRxiv 2017; doi:10.1101/207092.
- 343
344 14. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods
345 2012;9:357-359.
- 346
347 15. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella
348 K, Altshuler D, Gabriel S, Daly M *et al.* The Genome Analysis Toolkit: a MapReduce
349 framework for analyzing next-generation DNA sequencing data. Genome Res
350 2010;20:1297-1303.
- 351
352 16. Farek J, Hughes D, Mansfield A, Krasheninina O, Nasser W, Sedlazeck FJ, Khan Z,
353 Venner E, Metcalf G, Boerwinkle E *et al.* xAtlas: Scalable small variant calling across
354 heterogeneous next-generation sequencing experiments. bioRxiv 2018;
355 doi:10.1101/295071.
- 356
357 17. Lindner MS, Strauch B, Schulze JM, Tausch SH, Dabrowski PW, Nitsche A, Renard
358 BY. HiLive: real-time mapping of illumina reads while sequencing. Bioinformatics
359 2017;33:917-919.
- 360
361 18. McGann P, Bunin JL, Snesrud E, Singh S, Maybank R, Ong AC, Kwak YI, Seronello
362 S, Clifford RJ, Hinkle M *et al.* Real time application of whole genome sequencing for
363 outbreak investigation - What is an achievable turnaround time? Diagnostic
364 microbiology and infectious disease 2016;85:277-282.
- 365
366 19. Votintseva AA, Bradley P, Pankhurst L, Del Ojo Elias C, Loose M, Nilgiriwala K,
367 Chatterjee A, Smith EG, Sanderson N, Walker TM *et al.* Same-Day Diagnostic and
368 Surveillance Data for Tuberculosis via Whole-Genome Sequencing of Direct
369 Respiratory Samples. Journal of clinical microbiology 2017;55:1285-1298.

371 20. Greninger AL, Naccache SN, Federman S, Yu G, Mbala P, Bres V, Stryke D, Bouquet
372 J, Somasekar S, Linnen JM *et al.* Rapid metagenomic identification of viral pathogens
373 in clinical samples by real-time nanopore sequencing analysis. *Genome medicine*
374 2015;7:99.

375

376