

1 **Demonstrating the utility of flexible**
2 **sequence queries against indexed short reads**
3 **with FlexTyper**

4
5
6
7 **Phillip A. Richmond^{1*}, Alice M. Kaye^{1*}, Godfrain Jacques Kounkou¹, Tamar V. Av-**
8 **Shalom¹, Wyeth W. Wasserman¹**

9
10 ¹Centre for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute,
11 University of British Columbia, Vancouver, BC, Canada

12 *These authors contributed equally.

13
14 **Corresponding Author**

15 Wyeth W. Wasserman <wyeth@cmmt.ubc.ca>

16 Abstract

17 Across the life sciences, processing next generation sequencing data commonly relies
18 upon a computationally expensive process where reads are mapped onto a reference sequence.
19 Prior to such processing, however, there is a vast amount of information that can be ascertained
20 from the reads, potentially obviating the need for processing, or allowing optimized mapping
21 approaches to be deployed. Here, we present a method termed FlexTyper which facilitates a
22 “reverse mapping” approach in which high throughput sequence queries, in the form of k-mer
23 searches, are run against indexed short-read datasets in order to extract useful information. This
24 reverse mapping approach enables the rapid counting of target sequences of interest. We
25 demonstrate FlexTyper’s utility for recovering depth of coverage, and accurate genotyping of
26 SNP sites across the human genome. We show that genotyping unmapped reads can correctly
27 inform a sample’s population, sex, and relatedness in a family setting. Detection of pathogen
28 sequences within RNA-seq data was sensitive and accurate, performing comparably to existing
29 methods, but with increased flexibility. We present two examples of ways in which this
30 flexibility allows the analysis of genome features not well-represented in a linear reference. First,
31 we analyze contigs from African genome sequencing studies, showing how they distribute across
32 families from three distinct populations. Second, we show how gene-marking k-mers for the
33 killer immune receptor locus allow allele detection in a region that is challenging for standard
34 read mapping pipelines. The future adoption of the reverse mapping approach represented by
35 FlexTyper will be enabled by more efficient methods for FM-index generation and biology-
36 informed collections of reference queries. In the long-term, selection of population-specific
37 references or weighting of edges in pan-population reference genome graphs will be possible
38 using the FlexTyper approach. FlexTyper is available at
39 <https://github.com/wassermanlab/OpenFlexTyper>.

40

41 Author Summary

42 In the past 15 years, next generation sequencing technology has revolutionized our
43 capacity to process and analyze DNA sequencing data. From agriculture to medicine, this
44 technology is enabling a deeper understanding of the blueprint of life. Next generation
45 sequencing data is composed of short sequences of DNA, referred to as “reads”, which are often
46 shorter than 200 base pairs making them many orders of magnitude smaller than the entirety of a
47 human genome. Gaining insights from this data has typically leveraged a reference-guided
48 mapping approach, where the reads are aligned to a reference genome and then post-processed to
49 gain actionable information such as presence or absence of genomic sequence, or variation
50 between the reference genome and the sequenced sample. Many experts in the field of genomics
51 have concluded that selecting a single, linear reference genome for mapping reads against is
52 limiting, and several current research endeavors are focused on exploring options for improved
53 analysis methods to unlock the full utility of sequencing data. Among these improvements are
54 the usage of sex-matched genomes, population-specific reference genomes, and emergent graph-
55 based reference pan-genomes. However, advanced methods that use raw DNA sequencing data
56 to inform the choice of reference genome and guide the alignment of reads to enriched reference
57 genomes are needed. Here we develop a method termed FlexTyper, which creates a searchable
58 index of the short read data and enables flexible, user-guided queries to provide valuable insights
59 without the need for reference-guided mapping. We demonstrate the utility of our method by
60 identifying sample ancestry and sex in human whole genome sequencing data, detecting viral
61 pathogen reads in RNA-seq data, African-enriched genome regions absent from the global
62 reference, and HLA alleles that are complex to discern using standard read mapping. We
63 anticipate early adoption of FlexTyper within analysis pipelines as a pre-mapping component,
64 and further envision the bioinformatics and genomics community will leverage the tool for
65 creative uses of sequence queries from unmapped data.

66

67 Introduction

68 Short-read DNA sequencing enables diverse molecular investigations across life science
69 applications spanning from medicine to agriculture. Obtaining useful information from a data set
70 of raw reads (short pieces of DNA read outs from the DNA sequencer) typically involves
71 performing either *de novo* assembly, or mapping the read sequences against one or more
72 reference genomes. Whether the focus is on quantification (e.g. observed gene expression in
73 RNA sequencing data), or identifying sequence differences between a sample and a reference
74 genome (e.g. genotyping), the availability of a curated reference genome has led to a large
75 proportion of data analysis pipelines leveraging an indexed reference genome to perform
76 efficient read mapping as a primary analysis component.

77 Recently a plethora of large-scale, population-specific sequencing projects have
78 highlighted the numerous deficiencies and biases inherent to a single haploid reference (Ballouz,
79 Dobin et al. 2019, Yang, Lee et al. 2019). Examples include the large amount of structural
80 variation that exists between populations (Feuk, Carson et al. 2006, MacDonald, Ziman et al.
81 2014, Levy-Sakin, Pastor et al. 2019), the identification of unique sequences missing from the
82 current reference genome (Sherman, Forman et al. 2019), and population specific difference in
83 common genetic variants. Static linear reference genomes which do not capture these large
84 differences between populations impose challenges for accurate genotyping, with implications in
85 medicine and association studies (Ballouz, Dobin et al. 2019, Yang, Lee et al. 2019). Global
86 efforts to enrich the linear reference genome have led to the development of graph based
87 representations of pan-genomes, for a comprehensive review of current approaches see (Eizenga,
88 Novak et al. 2020, Sherman and Salzberg 2020). As highlighted in an earlier review by (Paten,
89 Novak et al. 2017), a key challenge in the future will be to determine the most appropriate
90 reference genome(s), or path(s) through a graph pan-genome, to maximize genotyping
91 performance. Knowledge regarding the genotypes of single nucleotide polymorphisms (SNPs) or
92 other makers present in a read data set can be used to guide the choice of reference.

93 Currently, the approach of identifying SNP genotypes across the genome primarily
94 involves computationally expensive reference-based read mapping and variant calling strategies
95 (Nielsen, Paul et al. 2011). Recently published tools have highlighted the expanse of information
96 that can be obtained from short read datasets. Inferring ancestry from specific, population-

97 discriminating SNPs can be performed rapidly with Peddy, which uses fewer than 25,000 SNPs
98 to identify ancestry through principal component analysis (Pedersen and Quinlan 2017).
99 Somalier (Pedersen, Bhetariya et al.) avoids the final stage of variant calling and evaluates
100 relatedness in aligned sequencing datasets. However, the accuracy of both of these tools is
101 affected by the underlying alignment. Previous work has shown that it is possible to genotype
102 predefined SNPs from unmapped sequence data, circumventing the read mapping and variant
103 calling process (Shajii, Yorukoglu et al. 2016, Dolle, Liu et al. 2017, Sun and Medvedev 2019).
104 Some approaches focus on k-mer (short sequences of length k) hashing and matching to
105 predefined target k-mers to perform genotyping of known SNPs, as demonstrated in the VarGeno
106 and LAVA frameworks (Shajii, Yorukoglu et al. 2016, Sun and Medvedev 2019). These
107 approaches are fast, but rely upon indexes of k-mers extracted from the reference genome and
108 SNP databases, thus reducing their flexibility for k-mers of different length and source. As we
109 move into the era of precision medicine, avoiding inherent reference bias is crucial in obtaining
110 accurate results. A separate approach is taken by Dolle et al., wherein the entire 1000 Genomes
111 dataset is compressed into an FM-index and queried with k-mers spanning polymorphic sites,
112 thus demonstrating the utility of scanning unmapped reads for predefined k-mers of interest. The
113 “reverse mapping” highlighted in their approach was applied to aggregated data, but the concept
114 can be extended to the analysis of individual genomes if implemented in a flexible way for
115 diverse types of queries.

116 The reverse mapping approach switches the focus onto querying for sequences of interest
117 within a read set, rather than a reference genome or database. This approach allows for a flexible
118 exploration of the information content of the reads by allowing the read set to be queried for
119 different parameters and across diverse sets of informative sequences. One example of this is
120 within RNA sequencing (RNA-seq), where analysis of cancer RNA-seq datasets can reveal the
121 presence of viral pathogens within patient data (Klijn, Durinck et al. 2015). Several tools have
122 been developed to specifically detect these viral pathogens from sequencing data including
123 viGEN (Bhuvaneshwar, Song et al. 2018) and VirTect (Xia, Liu et al. 2019). However, as with
124 the tools mentioned earlier, they are hampered by a mapping procedure which first maps against
125 the human reference genome and then subsequently maps against viral genome collections.
126 Other methods, such as Centrifuge (Kim, Song et al. 2016) and Kraken2 (Wood, Lu et al. 2019),
127 rely upon probabilistic or exact k-mer searches against large viral and bacterial databases. Both

128 of these methods are powerful, but lack flexibility and rely upon phylogenetic relationships
129 between target sequences. Specifically, they require the index for a search database to be
130 recreated for different k-mer lengths or when additional target sequences are added to the
131 database. Nevertheless, these tools are broadly used and thus serve as good comparators for
132 efficacy, as they have both been demonstrated to have utility in detecting viral pathogens within
133 cancer RNA-seq datasets by examining k-mer content.
134 (<https://www.sevenbridges.com/centrifuge/>).

135 Combining the current drive to decrease our reliance upon linear reference genomes, and
136 the wealth of demonstrated utility of reverse mapping approaches, we developed FlexTyper.
137 FlexTyper is a computational framework which enables the flexible indexing and searching of
138 raw next generation sequencing reads. We show example usage scenarios for FlexTyper by
139 demonstrating the high accuracy of reference-free genotyping of SNPs in single samples, and the
140 ability to identify foreign pathogen sequences within short-read datasets. We further explore the
141 utility of FlexTyper within challenging genomic regions hampered by hyper-variability, and test
142 its capacity to detect population-specific sequences missing from the reference genome. We hope
143 the flexibility afforded by the framework underpinning FlexTyper will fuel the emerging trend
144 away from the necessity for a static reference genome that currently lies at the heart of the
145 majority of genomic analysis tools.

146 Design and Implementation

147 Overview of FlexTyper

148 Usage can be broken down into three steps: 1) query generation, 2) indexing the raw
149 reads, and 3) querying against the FM-index (Figure 1). For query generation, we allow for both
150 custom user query generation, as well as pre-constructed queries from useful databases, such as
151 CytoScanHD array probe queries. Custom queries designed to capture genomic loci can be
152 generated by pairing a user-provided VCF (format v4.3) with a reference genome fasta file. For
153 the capture of potential pathogen sequences, we also allow query generation from one or more
154 fasta files. The files produced from query generation are used as input for subsequent index
155 query operations. The second step is the production of an FM-index from a set of short-read

156 sequences in fastq, gzipped fastq or plain text format. There is the option to include the reverse
157 complement of the reads within the index, however this increases the compute burden of
158 indexing, without the same reduction in search times. The read set is concatenated using a
159 sentinel character and passed as a single string into the indexing algorithm. The third step is the
160 core FlexTyper search algorithm which takes the query input file, generates search k-mers, and
161 scans the FM-index for matches. This step creates an output with matching format to the input
162 query file, with appended counts of matching reads for each query. A detailed breakdown of
163 these three components is described below.

164
165

166 **Fig 1 - Overview of FlexTyper**

167 FlexTyper has three primary components: query generation, read indexing, and searching against
168 the FM-index. Query generation includes the capacity to translate VCF files into query files
169 given a reference genome file (e.g. Genome Fasta), or to directly create queries from fasta
170 sequences including pathogen genome sequences. Modules VCF2Query.py and Fasta2Query.py
171 facilitate this process. The second component involves creating an FM-index of the raw reads,
172 after optional preprocessing steps. The third component searches the queries against the FM-
173 index to produce output files with counts of query sequences within the query files.

174

175 **Query generation**

176 FlexTyper supports flexible query generation giving users the capacity to query for any
177 target sequence or allele within their read dataset. Query files can be generated from an input
178 fasta and VCF file (VCF2Query.py), or directly from a fasta file (Fasta2Query.py). Potentially
179 useful queries, including those presented here, are provided online and include all sites from the
180 CytoScanHD chromosomal microarray, and ancestry discriminating sites (Pedersen and Quinlan
181 2017). These predefined query sets are available through git-lfs in the online FlexTyper Github
182 repository (<https://github.com/wassermanlab/OpenFlexTyper>). If users wish to directly query a
183 short-read dataset with a set of predetermined k-mers, we provide a separate function, ksearch,
184 that will directly search for k-mers from a given text file within an indexed read set.

185 **FM-index creation**

186 Generating the FM-index from short-read sequencing datafiles is performed in two steps:
187 preprocessing and indexing. The focus of our work is not on the algorithms used to construct the

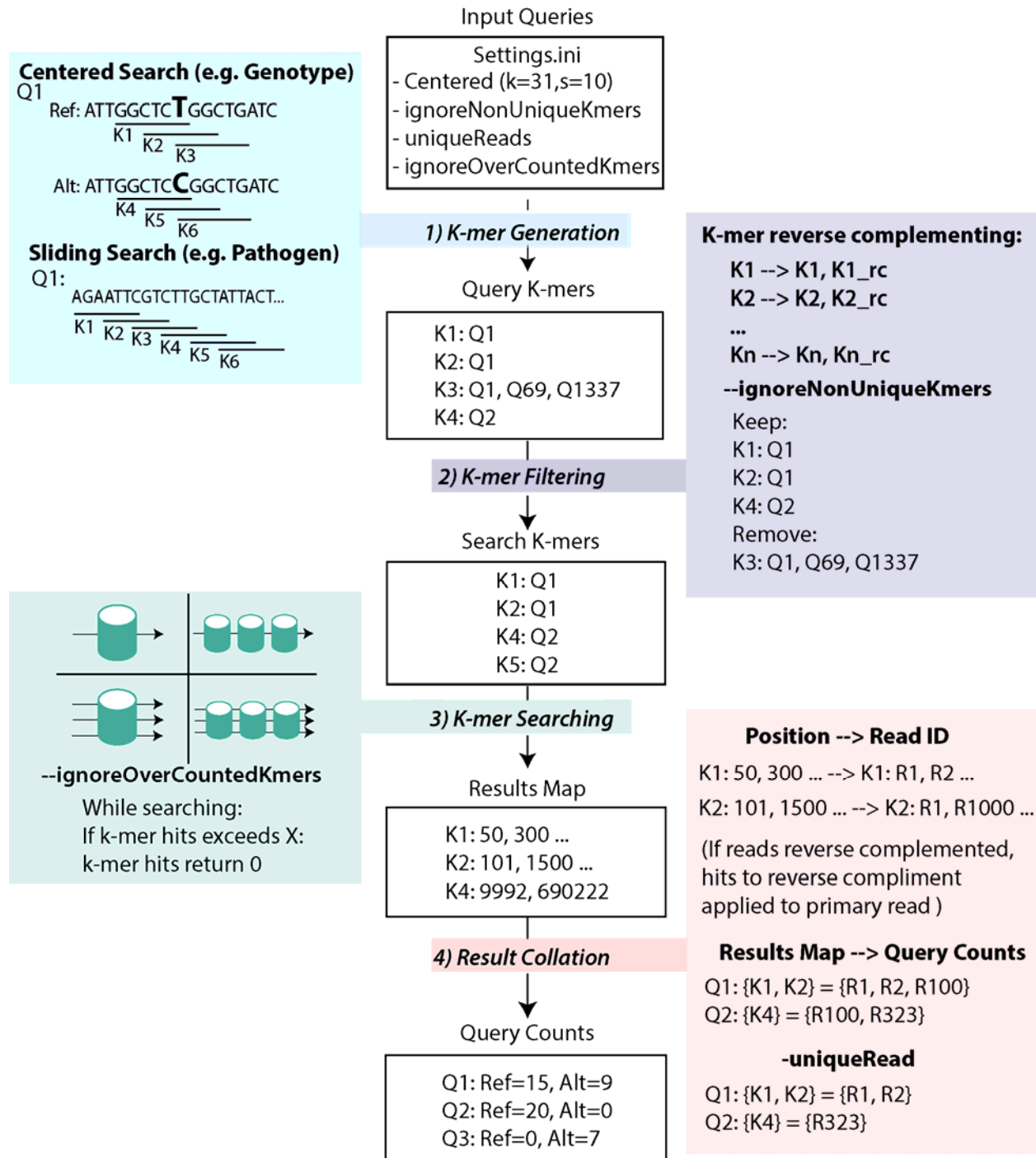
188 FM-index, and hence we use two existing utilities to generate a compatible FM-index for
189 FlexTyper. The toolkit Seqtk is used for reformatting the input read files by removing quality
190 scores and non-sequence information to create a sequence-only fasta format. The output fasta file
191 is processed using the SDSL-Lite library to generate the FM-index. SDSL builds a suffix array
192 that is used to generate the BWT of the input string, which is then compressed using a wavelet
193 tree and subsampled. The resulting compressed suffix array is streamed to a binary index file. As
194 the memory requirements for indexing large files can be burdensome, we support an option to
195 split the input file and index each chunk of reads independently. Downstream search operations
196 support the use of multiple indexes.

197

198 Query against FM-index

199 Querying the FM-index for user selected sequences can be conceptually divided into four
200 steps: 1) k-mer generation; 2) k-mer filtering; 3) k-mer searching; and 4) result collation (Figure
201 2). There are two primary methods of k-mer generation for a query; a centered search where the
202 middle position of the query is included in all k-mers, and a sliding search which starts at one
203 end of the query and uses a sliding window approach to generate the k-mers (Figure 2). Centered
204 search can be used for genotyping or estimating coverage over a single position, and the sliding
205 search can be used to count reads which match to any part of a query sequence. All parameters
206 for the search are specified in the settings.ini file, with a small number of key parameters able to
207 be overridden directly from the command line. After filtration, the k-mers are searched for within
208 the FM-index using C++ multithreading and asynchronous programming, using either a single
209 thread on a single index, multiple threads on a single index, a single thread on multiple indexes,
210 or multiple threads on multiple indexes (Figure 2). Importantly, asynchronous programming
211 allows the number of threads used during searching to be increased beyond the number of
212 available CPUs. The output from this search process is a collated results map containing the
213 positions of each k-mer within the FM-index. These positions are translated to read IDs, and
214 finally collapsed into query counts using the k-mer-query map. Importantly, if multiple k-mers
215 from the same query hit the same read, they are recorded as a single count at the query level. For
216 cases of multiple indexes being searched in parallel, the k-mer searching is performed
217 independently for each index, and then the search results from all indexes are merged and

218 reconciled to produce a final query count table. For a detailed explanation of the effects of key
 219 parameters, please see Supplementary Information, and our documentation on Github pages
 220 (<https://wassermanlab.github.io/OpenFlexTyper/>)
 221



222

223 **Figure 2 - Query Search Workflow**

224 Workflow for query search against the FM-index, starting with input queries and settings defined
225 in Settings.ini file. In this figure, the example shows a centered search with
226 ignoreNonUniqueKmers enabled. 1) K-mer generation has two modes centered search and
227 sliding search. For a centered search, the position of interest lies in the middle of the query, and
228 k-mers are designed to overlap that central position with defined length (k) and step (s). 2) If the
229 ignore-duplicates option is set, k-mers collated from the query set are filtered to remove any k-
230 mers which were found in multiple query sequences. 3) The filtered k-mers are then searched for
231 within a single FM-index (left two panels) or multiple indexes (right two panels) of the read set.
232 This can be done using single (top two panels) or multiple (bottom two panels) threads. 4) The
233 results corresponding to a position within the FM-index are then translated back into reads, with
234 hits on reverse complement reads assigned to the primary read, and collapsed into a set for each
235 query. The final counts are reported per query.
236

237 Post-processing of results into downstream formats

238 The output tables from the search process for genotyping can be translated into useful formats
239 for downstream analysis using the fmformatter scripts
240 (<https://github.com/wassermanlab/OpenFlexTyper/tree/master/fmformatter>). Currently, there is
241 the capacity to output genotype calls in VCF, 23andMe, or Ancestry.com format. Genotype calls
242 are derived here using a basic approach which assigns genotypes given a minimum read count
243 parameter as follows:

244 Alt < minCount && Ref > minCount: Homozygous reference, 0/0

245 Alt > minCount && Ref > minCount: Heterozygous alternate, 0/1

246 Alt > minCount && Ref < minCount: Homozygous alternate, 1/1

247 For searches which do not pertain to genotyping, the output tab-separated files can be used as
248 count tables for observed query sequences.

249

250

251

252

253 Results

254 Observations about FlexTyper system requirements and performance

255 The generation of a full text index of the reads is a key step and we are able to generate
256 indexes of human whole genome sequencing reads with ~800 million reads utilizing less than
257 150Gb of RAM on a single compute node within a higher performance compute (HPC) cluster.
258 (Supplemental Table S1). Although read indexing is slower than a traditional alignment, sorting,
259 and variant calling pipeline, FlexTyper can index whole genome sequencing samples in roughly
260 24 hours (Supplemental Table S1). The flexibility for creating sub-indexes allows the user to
261 adjust parameters to fit their system, accommodating most modern HPC architectures. In
262 comparison to the tools that utilize prebuilt indexes, such as BWA-MEM, or generate
263 probabilistic indexes, such as Kraken2, FlexTyper is significantly slower, however, the non-
264 approximate full text implementation allows the read set to be queried for diverse sequences,
265 across the full parameter space, without reindexing unlike Kraken2. The index of a high depth
266 paired-end (2x250bp) WGS read set for sample HG002 uses only 155 Gb RAM, and with no
267 information loss in the read sequences, it is not necessary to retain the original fastq files once
268 indexed. While FlexTyper does remove quality scores, the tuning of k-mer length and step size
269 allows counting of a read even in the presence of errors, and filtering the fastq file to remove
270 low-quality/high-error reads can be done prior to read indexing. The complex interplay of the
271 different search parameters makes generalized performance statements challenging, so to inform
272 the user of FlexTyper's use-case performance we provide runtime metrics and read recovery for
273 a variety of different search settings and scenarios in the following sections.

274

275 Testing for the presence of pathogen sequences in RNA-seq

276 To demonstrate the capacity of FlexTyper to detect pathogens from RNA-sequencing
277 data, we generated synthetic reads from four relevant viral genomes including Epstein-Barr virus
278 (EBV), Human Immunodeficiency virus type 1 (HIV-1), and two Human Papilloma virus strains
279 68b (HPV FR751039) and 70 (HPV U21941) (Supplemental Methods). We first examined the
280 impact of various FlexTyper parameters on the recovery rate of pure, simulated read sets for each

281 of the four viruses and one human blood RNA-seq dataset from the Genome England project
282 (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-6523/samples/>) (Table S2).
283 Importantly, varying the parameters k (length of the k-mer search substring) and s (step-size)
284 change the specificity and sensitivity of read recovery. When k is set to 15 (a short k-mer), there
285 were roughly 1 million off-target hits to the viral genomes for the pure human RNA-seq file, and
286 the increased search space leads to increased run times (Table S2). Next, we explored the effect
287 of two uniqueness settings, the first for k-mers, (ignoreNonUniqueKmers), where a given k-mer
288 cannot map across multiple queries, and secondly for reads, (uniqueReads), where a read cannot
289 be counted across multiple queries. We provide both parameters to the users, as there may be
290 instances where reads are allowed to be counted across queries, but the k-mers must be
291 independent. By exploring these parameters, we show that all simulated reads can be recovered
292 with parameters of $k=30$ and $s=5$, and off-target assignment can be controlled with the
293 uniqueness settings. For more details about the parameter explanations or the impact of
294 parameterizations on read recovery, see Supplemental Information and Table S2.

295 Next, we simulated mock patients infected by the four viruses to examine the detection
296 capacity of FlexTyper with respect to two established methods, Centrifuge and Kraken2.
297 Simulated reads from each of the four viruses were spiked-in at various concentrations (read
298 counts) within five different human RNA-seq datasets from the Genome England project
299 (<https://www.personalgenomes.org.uk/data/>) (Supplemental Methods). Using the optimized
300 parameters derived above ($k=30$, $s=5$, uniqueRead, and uniqueK), we are able to detect each
301 virus in the patient samples even at low concentrations (Figure 3). We also searched these mock
302 patient samples using an expanded set of parameters, and see the expected changes in sensitivity
303 as the uniqueness and k parameters are altered (Table S3). For more direct comparison to
304 Centrifuge and Kraken2 we ran FlexTyper in the paired-read mode, and show that we are able to
305 detect the pathogen sequences with high sensitivity and specificity (Figure 3). This is true even
306 for a sample (Patient 5) which had a 1x concentration—roughly 50-1150 reads depending on
307 genome size—of each viral genome spiked in. For Patient 4, where the level of EBV spiked in is
308 over 1 million reads, the undercounting stems from our limit on maximum occurrence, which we
309 set to 2000, and adjusting to a maximum occurrence of 10,000 increases the count to 1,144,990.
310 For the HPV viruses (U21941 and FR751039), we are able to detect the levels of both strains,
311 and observe a slight over-counting of FR751039 in Patient 2, possibly from a high count of

312 spiked in U21941 reads. We compared our approach with Centrifuge and Kraken2, which match
313 reads based on k-mers mapped against a comprehensive indexed viral and bacterial database, and
314 tabulate matches at the read pair level. Centrifuge works with unmapped short-read sequencing
315 data by performing read-length (k=150) k-mer searches against a database of viral and bacterial
316 genomes, and hence was the least sensitive method due to the limitations of full length k-mer
317 queries (Kim, Song et al. 2016). Kraken2, which uses a minimizer for approximate matching,
318 can search for shorter k-mers (default k=31), leading to increased sensitivity over the Centrifuge
319 method. Both Centrifuge and Kraken2 achieve accurate results for the HIV-1 and EBV samples,
320 but were only able to detect 5-10% of the reads for the two HPV samples, U21941 and
321 FR751039, even when aggregating at the family level (Papillomaviridae) (Figure 3). Initially, we
322 hypothesized that this was due to off-target mapping to another genome, perhaps the human
323 genome, within their comprehensive database. However, after testing a set of 5200 pure U21941
324 or FR751039 reads with Kraken2, we were only able to recover 136 reads and 279 reads
325 respectively, even when considering reads assigned to the viral kingdom level (Supplemental
326 Methods). This limitation can be overcome with FlexTyper, which enables the user to define the
327 relevant pathogens to search for, along with the ability to repeat searches across different k-mer
328 lengths, without the need for re-indexing a complex bacterial or viral database. While Kraken2
329 and Centrifuge are powerful and comprehensive metagenomic classifiers, which allow for
330 increased breadth and classification across a phylogenetic tree, there may be cases where specific
331 pathogen queries of interest require high sensitivity. We believe FlexTyper can serve as an
332 option in these scenarios.
333



334

335 **Figure 3 - Mixed Viral Analysis**

336 Detection of pathogen sequences in five synthetic patient RNA-seq datasets (Patient 1-5; rows),
337 each with different levels of spiked-in viruses (EBV, HIV-1, U21941, and FR751039; columns),
338 expected values shown as black vertical bars. As Centrifuge and Kraken2 are unable to delineate
339 between the two HPV substrains (U21941 and FR751039), a combined count at the HPV level is
340 tabulated.

341

342 Genomic coverage and genotype detection within human WGS data

343 Knowing whether a given k-mer is present or absent from a human WGS datafile (in this
344 instance Illumina short-read, paired-end data) can have utility for estimating the depth of
345 coverage for a target region and genotyping SNPs. FlexTyper has the capacity to compute depth
346 of coverage or genotype SNPs from WGS data for both predefined and user-supplied loci. We
347 demonstrated this capacity for genomic sites using the probe sequences from the CytoScanHD

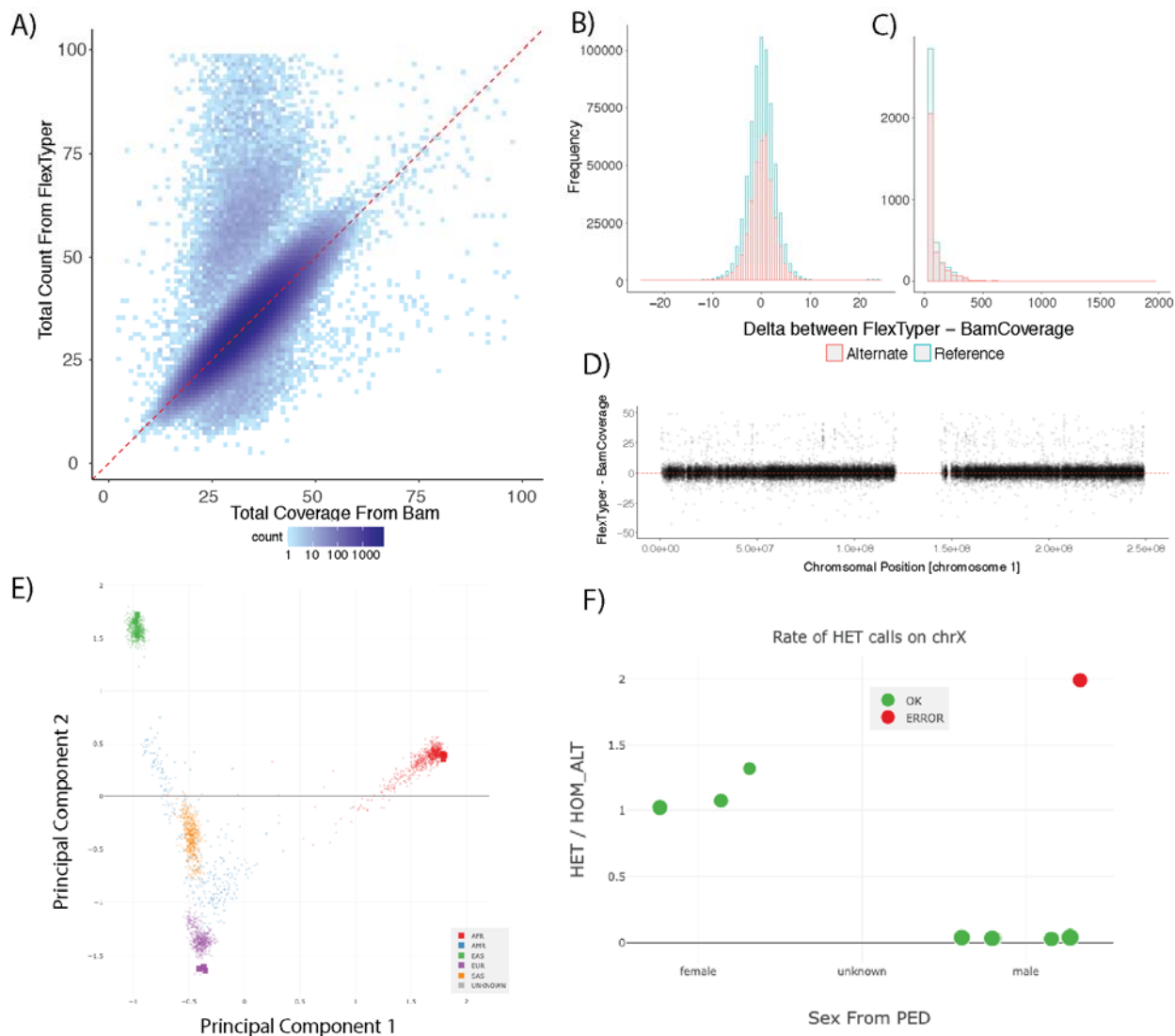
348 microarray, as well as a subset of previously collated population discriminating SNPs (Pedersen
349 and Quinlan 2017). Using these loci, we created query files with a reference and alternate query
350 sequence centered on the biallelic site (Supplemental Methods).

351 We first sought to test the read recovery capacity of FlexTyper compared to an alignment
352 based method which we call BamCoverage. The BamCoverage method involves mapping the
353 reads to the reference genome, and then extracting per-base read coverage over a specific
354 reference coordinate. BamCoverage utilizes the pysam package to extract read pileup over
355 positions defined by the FlexTyper input query file (Supplemental Methods). Using the
356 CytoScanHD SNP set, we found a high concordance between the read counts from FlexTyper
357 and the depth of coverage from aligned reads (Figure 4A). FlexTyper was run with parameters of
358 $k=31$, $s=10$, max occurrence 200, and the requirement of unique k-mers between queries. The
359 vast majority, 98.33% (784,297/797,653), of sites differed by less than 10 between FlexTyper
360 and BamCoverage, with a Spearman correlation of 0.86 (Figure 4B). The discrepancy in counts
361 is similar for both reference and alternate alleles, which is important since most genotyping
362 models assume relative contributions of observed alleles for genotype calling. There were 10,397
363 sites with a delta, ($\Delta = \text{FlexTyper} - \text{BamCoverage}$), greater than 10, and 1,469 sites with a delta
364 greater than 100 (Figure 4C). We manually investigated a few of these sites which were
365 overcounted by FlexTyper by more than 100 and found that they are being overcounted due to k-
366 mers mapping to multiple possible locations. Comparing these over-counted hits with delta
367 greater than 100 to previously defined repeat regions shows that 1444/1469 or 98.3% of the
368 overcounted sites overlapped with predefined repeats (Trost, Walker et al. 2018). The uniqueness
369 of k-mers is important for accurate read counting, thus it is recommended to filter query
370 sequences within such regions when using FlexTyper for genotyping or depth profiling. Lastly,
371 by examining the recovery of reads across the chromosome between FlexTyper and the
372 BamCoverage approach, we observe uniform recovery across the breadth of the chromosome
373 (Figure 4D). This is important for copy number variant calling applications, as they rely upon
374 contiguous readouts of genomic sequence coverage.

375 Next, we investigated whether FlexTyper can accurately recover genotypes at the SNP
376 sites profiled from the chromosomal microarray. The genotyping approach we use leverages a
377 minimum count from the reference and alternate allele to assign heterozygous, homozygous
378 alternate, or homozygous reference genotypes (Supplemental Methods). We applied this basic

379 genotyping algorithm to both FlexTyper and BamCoverage counts to produce a VCF file. These
380 genotypes were compared to an alternate pipeline which uses reference-based mapping and
381 sophisticated variant calling using BWA-MEM and DeepVariant (Poplin, Chang et al. 2018)
382 (Supplemental Methods). For the CytoScanHD microarray, all three methods report the same
383 genotype for 99.4% (792,805/797,653) of the SNP sites. We further investigated the discordant
384 genotypes to see if we can explain why there is a disagreement between FlexTyper and the other
385 two methods. First, we compared these genotypes to the repeats defined above and see that
386 77.6% (2,980/3,838) of the discordant genotypes overlap with predefined repeat loci.
387 Additionally, we implemented a flag for offending k-mers, to signal when a k-mer was non-
388 unique or surpassed the max occurrence (maxOcc) parameter. There were a total of 685
389 genotypes with offending k-mers, of which 355 (8.6% of the 3,838) overlap with the discordant
390 genotypes unique to FlexTyper. We further demonstrate the accuracy of FlexTyper-derived
391 genotypes by indexing nine WGS samples from the Polaris project representing diverse
392 populations including three African, three Southeast Asian, and three European individuals
393 (Chen, Krusche et al.). After indexing, we queried the samples for population discriminating
394 sites and then genotyped the output table to produce a VCF file. The output VCFs were then used
395 within the Peddy tool, and a principal component analysis was performed to predict the ancestry
396 of the samples (Pedersen and Quinlan 2017). In all nine cases the population was correctly
397 determined, as well as the relatedness inference for the three trios (Figure 4E, Figure S1).
398 Interestingly, we observed a discrepancy between the listed sex for the child of the European
399 trio, individual HG01683, and the inferred sex from FlexTyper and Peddy (Figure 4F). We
400 followed up on this observation and revealed that the individual is not an XY male, but rather an
401 XXY individual, and communication was made that resulted in the relabeling of the individual
402 within the online repository. The analysis time for extracting the queries from the indexed reads
403 for all WGS analysis can be found in Table S4, highlighting that especially for informative
404 subsets of queries, such as population discriminating sites, genotypes can be recovered
405 accurately and quickly (~10-15 minutes). Taken together, FlexTyper has the capacity to provide
406 accurate counts of observed reads matching two alleles over informative SNP sites, with relevant
407 utilities such as copy number estimation, sample identification, ancestry typing, and sex
408 identification.

409



410

411 **Figure 4 - WGS Genotyping using FlexTyper**

412 A) FlexTyper read count compared to the total coverage from BAM file over SNP sites
 413 represented on the CytoScanHD microarray. B) Histogram showing the delta, ($\Delta = \text{FlexTyper} -$
 414 BamCoverage), in read count for both the alternate (red) and reference (blue) alleles. C)
 415 Histogram of the same delta as B) but with an extended axis from 100-2000, showing the
 416 frequency of over-counting for sites using FlexTyper. D) Scatter plot showing the delta ($\Delta =$
 417 $\text{FlexTyper} - \text{BamCoverage}$) on the y-axis, plotted across chromosome 1 on the x-axis. E)
 418 Principal component analysis showing projection of FlexTyper-derived SNP genotypes from
 419 nine individuals of Asian (green), African (red) and European (purple) ancestry. Squares denote
 420 FlexTyper genotypes, points denote existing data from the 1000 Genomes project provided by
 421 Peddy. F) Sex-typing for these Polaris samples showing the ratio of heterozygous to homozygous
 422 sites on the X chromosome (y-axis) for individuals for the defined sexes as male (right) and
 423 female (left). Each individual is labeled as green (correctly sex-labeled) or red (incorrectly
 424 labeled).

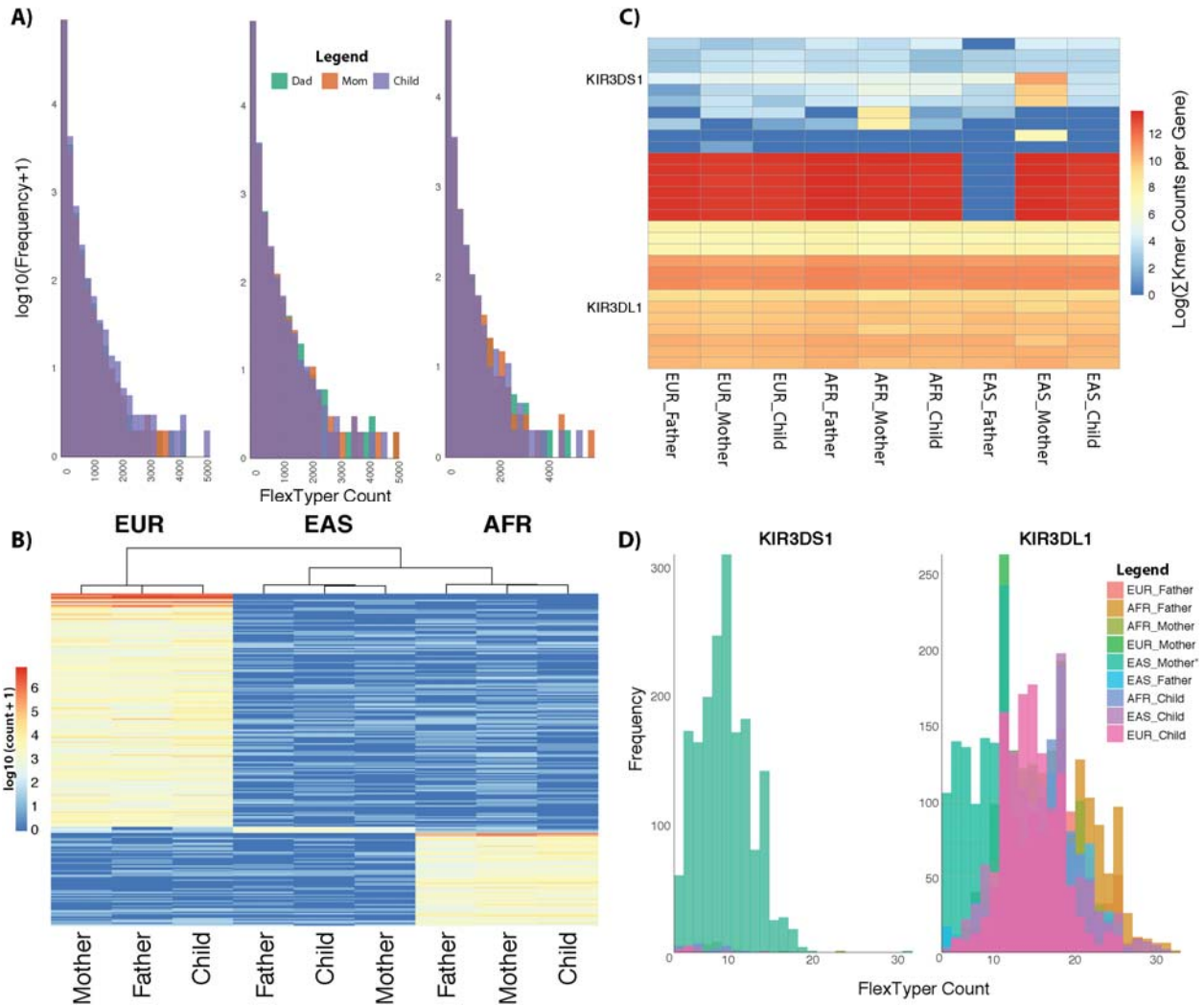
425 Exploring creative uses of FlexTyper

426 To demonstrate the flexible utility of our k-mer-based searching method, we explored
427 regions of the genome which are challenging for read mapping and downstream analysis when
428 represented within haploid, linear reference genomes. The two areas we chose to focus on
429 include contigs derived from a population but not present in the reference genome, and
430 hypervariable and homologous regions, where linear representations are known to perform
431 poorly.

432 The contigs we chose to process include the “non-reference” contigs from a recent
433 African pan-genome publication (Sherman, Forman et al. 2019). These contigs, which were
434 assembled from non-mapped reads, collectively contain nearly 300 megabases of DNA,
435 represented by ~125,000 contigs. We created queries of these contigs, and then searched them
436 using the nine samples from our ancestry WGS experiment using parameters of $k=50$, $s=5$, and
437 `uniqueRead=true`. Unexpectedly, when we queried the African contigs across the family trio
438 samples from three populations (EAS, EUR, and AFR), we observed similar contig coverage
439 across all individuals (Figure 5A). Next, we sought to identify discriminating contigs within the
440 ~125,000 non-reference contigs by filtering for those which consistently appear in one
441 population (>10 counts in mother, father, and child), but had low coverage in the other two
442 populations (<5 counts). Applying this to the three groups, we identified a set of discriminating
443 contigs (Figure 5B). The African trio had 151 unique contigs, the East Asian trio had 60 unique
444 contigs, and only four contigs were unique to the European trio. In total, the analysis indicates
445 that the African-derived contigs are widespread across populations. Two limitations of our
446 analysis include: 1) FlexTyper was run in unique mode, so reads mapping across highly similar
447 contigs are discounted, and 2) FlexTyper does not account for local genome context, so it is
448 possible that some of the contigs are unique not due to specific sequence, but due to their
449 placement in the genome (e.g. structural variants). This application highlights the potential of
450 FlexTyper in filtering and querying for contigs unique to a subpopulation.

451 Recent tailored approaches to genotyping challenging genomic regions, which are
452 difficult due to their hypervariability in the population and/or high sequence similarity between
453 homologous genes, utilize unique k-mer counts to distinguish between alleles present in a sample
454 (Roe and Kuang, Shen and Kidd 2020). As FlexTyper has the capacity to rapidly query k-mers
455 and generate unique k-mers across input queries, we decided to test FlexTyper’s utility in

456 distinguishing between samples for a locus known to be challenging: the killer-cell immune
457 receptor (KIR) locus. We downloaded a curated set of gene-distinguishing k-mers for this locus
458 which have been used, with the k-mer counting tool KMC3, to identify the presence-absence of
459 the 28 genes/alleles in the KIR locus (Roe and Kuang). Using the FlexTyper function ksearch,
460 we searched the nine WGS samples for this set of k-mers, and then tallied the k-mers with >3
461 counts per gene (Figure 5C). While we did not see many differences at the family level with this
462 set, we did observe an outlier sample: the mother in the East Asian trio. For the KIR3DS1 gene,
463 she had several high-counting k-mers which were absent in the other samples. KIR3DS1 is an
464 alternate haplotype of the KIR3DL1 gene in the canonical reference genome, and is represented
465 in the GRCh38 reference on an ALT contig (chr19_KI270887v1_alt). By plotting a histogram of
466 the counts for the nine individuals over both the KIR3DS1 and KIR3DL1 genes, we observed
467 that the mother of the East Asian trio is the only sample in this set with k-mer coverage over
468 KIR3DS1, and consequently has reduced coverage of the KIR3DL1 gene (Figure 5D,E). Taken
469 together, this suggests that the mother is heterozygous for the KIR3DL1 and KIR3DS1 alleles,
470 while the rest of the samples in this set are homozygous for the KIR3DL1 allele. This
471 observation is enabled by the careful selection of k-mers by Roe et al., and is a demonstration of
472 how FlexTyper can utilize user curated k-mers for genotyping within a challenging locus.
473



474

475 **Figure 5 - Explorative uses of FlexTyper**

476 Two examples of the creative uses of FlexTyper within challenging regions. A) Histograms for
 477 counts of the African contigs overlaid for father (green), mother (orange) and child (purple) split
 478 by population (left to right EUR, EAS, AFR). B) Heatmap in log-scale for population-specific
 479 contigs, clustered by sample similarity (columns) and contig count similarity (rows). C) Heatmap
 480 showing the \log_{10} transform of the sum of k-mer counts per gene, with genes as rows, and
 481 samples as columns. The two alleles, KIR3DS1 and KIR3DL1 are labeled as rows on the left
 482 side. D) Overlaid histograms for the 9 samples, showing the frequency of the FlexTyper k-mer
 483 count for the KIR3DS1 (left) and KIR3DL1 (right) alleles.

484 **Discussion**

485 Here we presented FlexTyper, a user-friendly tool which enables exploratory analysis of
 486 short read datasets without the need to perform reference guided alignment. Our framework
 487 allows the user to generate custom queries, or to directly search from a list of k-mers. This gives

488 the user complete flexibility to tailor the search inputs and parameters to the problem at hand.
489 We demonstrated three common applications: depth of coverage analysis, accurate SNP
490 genotyping, and sensitive detection of pathogen sequences. We then showcased the potential for
491 FlexTyper to extract useful information from complex, hypervariable, or non-reference genomic
492 sequences. FlexTyper was designed with user simplicity in mind, but without comprising the
493 breadth of potential applications, and hence the tool is available for the creative use of genomics
494 researchers.

495 The rapid and accurate recovery of read depth enables innovative usage of FlexTyper in
496 the space of copy number variant profiling. We demonstrated that we could reproduce the depth
497 of coverage of a genomic region without the need for reference-based mapping. As microarrays
498 begin to be replaced by genome sequencing assays, we envision that FlexTyper could be
499 extended to reproduce microarray-style outputs that are established in clinical labs. Further, we
500 show that when genomic queries with counts higher than the expectation arise, these events
501 correspond to repetitive genomic sequences. As such, FlexTyper may not only enable the
502 recovery of read depth in an accurate manner, but it can also inform the quality of a sequence
503 query as a “unique probe” for assessing genomic copy number.

504 The genotyping case study highlights how pre-alignment analysis of genome sequence
505 data can provide rapid insights into the properties of a sample. SNP genotyping was accurate
506 across the genome, allowing rapid identification of sample ancestry, sample relatedness in the
507 trio setting, and sample sex typing using Peddy (Pedersen and Quinlan 2017). Interestingly,
508 applying Peddy to the output of FlexTyper for open source trio data from the Polaris project
509 revealed a mislabeling of the sex for individual HG01683, which was reported and subsequently
510 amended in the online data repository ([https://github.com/Illumina/Polaris/wiki/HiSeqX-Kids-](https://github.com/Illumina/Polaris/wiki/HiSeqX-Kids-Cohort)
511 [Cohort](https://github.com/Illumina/Polaris/wiki/HiSeqX-Kids-Cohort)). Since ancestry and sex information can inform choices in downstream data processing,
512 identifying these discrepancies between labeled sex and inferred sex in a data-driven manner is a
513 critical step of pre-alignment informatics. For instance, mapping against the sample-matched sex
514 chromosomes has been shown to improve performance (Olney, Brotman et al. , Webster, Couse
515 et al. 2019). As such, using FlexTyper, in combination with Peddy, on diverse datasets prior to
516 reference-guided read alignment will lead to improved results from mapping-based pipelines.

517 There is increased recognition of the important of pathogen detection. In both cancer
518 profiling (Klijn, Durinck et al. 2015) and public health studies (Gardy, Loman et al. 2015), rapid

519 determination of the presence of pathogen sequences could obviate the need for full reference
520 mapping. Some existing tools designed for viral detection in sequencing data rely upon pre-
521 indexed databases of viral and bacterial sequences, sometimes including a phylogenetic
522 relationship between genomes within the index (Kim, Song et al. 2016, Wood, Lu et al. 2019,
523 Xia, Liu et al. 2019). Two approaches, Centrifuge and Kraken2, have been applied to cancer
524 genomes to confirm the presence of viral pathogens, including Human papilloma virus (HPV).
525 We demonstrated that our approach compares favorably to Centrifuge, with a more sensitive
526 detection level, due to the ability to search for k-mers shorter than the read length and the
527 advantage of fine-tuned control over the searchable database. Comparing FlexTyper to Kraken2,
528 which doesn't rely upon full read length queries, detection of the spiked-in pathogen sequences
529 was as good or better than Kraken2, with improved performance for detecting the HPV-derived
530 reads. Interestingly, both Kraken2 and Centrifuge had difficulties in detecting HPV reads, both
531 within mixed-virus and pure viral read sets. Here we only searched for viral pathogens of
532 interest, although other specific pathogen queries could be performed, such as the presence of
533 antibiotic resistance genes within a patient RNA-seq sample.

534 As the research community begins to move away from a single haploid reference towards
535 richer pan-genome representations, we anticipate that more diverse and creative uses for
536 FlexTyper's 'reverse mapping' approach will emerge. During our continued exploration of
537 FlexTyper's potential, we have identified a few possible applications. We focused on regions of
538 the genome which are challenging for traditional linear reference approaches, including a set of
539 sequences not present in the reference genome, and a highly polymorphic region with
540 homologous genes. Using the set of contigs assembled from the African Pan-genome project, we
541 applied FlexTyper and observed similar sequence coverage over these contigs across three
542 families from European, East Asian, and African ancestry. We further filtered this set of contigs
543 and identified a limited set of discriminating contigs, highlighting another relevant use case for
544 FlexTyper. Beyond non-reference contig searching, we explore the utility of FlexTyper in
545 genotyping the polymorphic and homologous genes within the KIR locus. We use a curated set
546 of k-mers from a work by Roe et al, and identify an alternate haplotype (namely KIR3DS1,
547 present in the ALT contigs of GRCh38) for the KIR3DL1 gene in one of the nine individuals.
548 This genotyping demonstration with a curated set of k-mers highlights the potential for
549 FlexTyper to be adopted by other specialized methods tailored for challenging genomic regions.

550 The full breadth of possible applications of FlexTyper and its reverse mapping approach
551 has yet to be discovered, but we have highlighted multiple potential avenues here. For WGS read
552 data sets, it is feasible to genotype complex structural variants by searching for sequences
553 overlapping breakpoints, such as those observed in a subpopulation, or events recurrently found
554 in cancer (Sudmant, Rausch et al. 2015, Li, Roberts et al. 2020). Within RNA-seq data, querying
555 for exon-exon splice junctions in a rapid manner can allow isoform quantification, as has been
556 previously demonstrated (Patro, Mount et al. 2014, Bray, Pimentel et al. 2016). Further, a recent
557 report showed the utility of k-mer-counting methods in resolving copy number variants within
558 paralogous loci and genes (Shen, Shen et al. 2020). Another group showed the advantage of
559 examining depth of coverage at specific sites across the paralogous genes in Spinal Muscular
560 Atrophy (Chen, Sanchis-Juan et al. 2020) As FlexTyper is well suited for specific sequence
561 recovery operations, scanning with preselected query sequences such as defined by these studies
562 can enable rapid detection (Chen, Sanchis-Juan et al. 2020). All of these proposed applications
563 help tackle challenges which are currently a burden for traditional reference-based mapping
564 approaches.

565 We focused this report on the expansive utility of querying indexed read sets for
566 interesting and informative sequences, but recognize that speed and computational resources are
567 an important consideration in the adoption of the method. One obvious, but transient, constraint
568 on the utility of FlexTyper is the ability to generate the FM-index of a read data set. Our
569 implementation utilizes the SDSL library, chosen for its stability, however as the FM-index is
570 critical to many aspects of genome scale analyses there have been strong efforts to develop more
571 efficient indexing algorithms. Recent methods have shown both dramatic increases in
572 construction speed either through induced suffix sorting (Kärkkäinen, Kempa et al. 2017) or
573 GPU-based construction algorithms (Chacón, Marco-Sola et al. 2015), and decreases in memory
574 requirements (Labeit, Shun et al. 2017, Chen, Li et al. 2018). Within our current framework, to
575 try and mitigate some of these issues, we built in methods to split the read set into smaller
576 chunks, each of which is indexed in serial. Although not currently implemented, it is clear that
577 this could be executed in parallel, if there is sufficient RAM available, as each index is generated
578 independently of other chunks. Furthermore, the nature of the reverse mapping approach holds
579 promise with massive parallelization approaches, including those involving GPU acceleration
580 (Hung, Hsu et al. 2018). Moving forward, accelerations to the FM-index generation and reverse

581 mapping approach will result in faster genomic analysis pipelines than is currently possible with
582 alignment-based methods.

583 Looking to the future, we see the k-mer-searching approach of FlexTyper as having great
584 utility when used in conjunction with emergent pan-genome and graph representations of the
585 reference genome (Kehr, Trappe et al. 2014, Kaye 2016, Paten, Novak et al. 2017). Whether
586 users seek to select a population specific reference graph as the basis for read mapping, or to
587 introduce Bayesian priors (edge weighting) within a pan-population reference graph, knowledge
588 of population markers spanning chromosomes will be required to inform the processes.
589 Furthermore, it is our expectation that pan-genome mapping methods will ultimately use full text
590 read-based indexes, to allow for data compression without loss of information or functionality,
591 while avoiding the plethora of issues facing approaches that index a pan-genomic representation
592 (Paten, Novak et al. 2014, Ghaffaari and Marschall 2019). As the reference structure is enriched
593 and algorithms for use with pan-genome graphs mature, approaches such as FlexTyper, which
594 enable the reverse mapping of informative sequences against a set of indexed reads, will be
595 instrumental in the initial steps of genome analysis pipelines.

596
597

598 References

- 599
600 Ballouz, S., A. Dobin and J. A. Gillis (2019). "Is it time to change the reference genome?"
601 Genome Biol. **20**(1): 159.
- 602 Bhuvaneshwar, K., L. Song, S. Madhavan and Y. Gusev (2018). "viGEN: An Open Source
603 Pipeline for the Detection and Quantification of Viral RNA in Human Tumors." Front.
604 Microbiol. **9**: 1172.
- 605 Bray, N. L., H. Pimentel, P. Melsted and L. Pachter (2016). "Near-optimal probabilistic RNA-
606 seq quantification." Nat. Biotechnol. **34**(5): 525-527.
- 607 Chacón, A., S. Marco-Sola, A. Espinosa, P. Ribeca and J. C. Moure (2015). "Boosting the
608 FM-Index on the GPU: Effective Techniques to Mitigate Random Memory Access." IEEE/ACM
609 Trans. Comput. Biol. Bioinform. **12**(5): 1048-1059.

- 610 Chen, N., Y. Li and Y. Lu (2018). A Memory-Efficient FM-Index Constructor for Next-
611 Generation Sequencing Applications on FPGAs. 2018 IEEE International Symposium on
612 Circuits and Systems (ISCAS): 1-4.
- 613 Chen, S., P. Krusche, E. Dolzhenko, R. M. Sherman, R. Petrovski, F. Schlesinger, M. Kirsche,
614 D. R. Bentley, M. C. Schatz, F. J. Sedlazeck and M. A. Eberle (2019). "Paragraph: a graph-based
615 structural variant genotyper for short-read sequence data." Genome Biol. **20**(1): 291.
- 616 Chen, X., A. Sanchis-Juan, C. E. French, A. J. Connell, I. Delon, Z. Kingsbury, A. Chawla, A.
617 L. Halpern, R. J. Taft, N. BioResource, D. R. Bentley, M. E. R. Butchbach, F. L. Raymond and
618 M. A. Eberle (2020). "Spinal muscular atrophy diagnosis and carrier screening from genome
619 sequencing data." Genet. Med.
- 620 Dolle, D. D., Z. Liu, M. Cotten, J. T. Simpson, Z. Iqbal, R. Durbin, S. A. McCarthy and T. M.
621 Keane (2017). "Using reference-free compressed data structures to analyze sequencing reads
622 from thousands of human genomes." Genome Res. **27**(2): 300-309.
- 623 Eizenga, J. M., A. M. Novak, J. A. Sibbesen, S. Heumos, A. Ghaffaari, G. Hickey, X. Chang,
624 J. D. Seaman, R. Rounthwaite, J. Ebler and Others (2020). "Pangenome Graphs." Annu. Rev.
625 Genomics Hum. Genet. **21**.
- 626 Feuk, L., A. R. Carson and S. W. Scherer (2006). "Structural variation in the human genome."
627 Nature Reviews Genetics **7**(2): 85-97.
- 628 Gardy, J., N. J. Loman and A. Rambaut (2015). "Real-time digital pathogen surveillance —
629 the time is now." Genome Biology **16**(1).
- 630 Ghaffaari, A. and T. Marschall (2019). "Fully-sensitive seed finding in sequence graphs using
631 a hybrid index." Bioinformatics **35**(14): i81-i89.
- 632 Hung, C.-L., T.-H. Hsu, H.-H. Wang and C.-Y. Lin (2018). "A GPU-based Bit-Parallel
633 Multiple Pattern Matching Algorithm." 2018 IEEE 20th International Conference on High
634 Performance Computing and Communications; IEEE 16th International Conference on Smart
635 City; IEEE 4th International Conference on Data Science and Systems (HPCC/SmartCity/DSS).
- 636 Kärkkäinen, J., D. Kempa, S. J. Puglisi and B. Zhukova (2017). "Engineering External
637 Memory Induced Suffix Sorting." 2017 Proceedings of the Nineteenth Workshop on Algorithm
638 Engineering and Experiments (ALENEX).
- 639 Kaye, A. (2016). Methods for the graphical representation of genomic sequence data. US
640 Patent. Uspto, University of British Columbia.
- 641 Kehr, B., K. Trappe, M. Holtgrewe and K. Reinert (2014). "Genome alignment with graph
642 data structures: a comparison." BMC Bioinformatics **15**(1): 99.
- 643 Kim, D., L. Song, F. P. Breitwieser and S. L. Salzberg (2016). "Centrifuge: rapid and
644 sensitive classification of metagenomic sequences." Genome Res. **26**(12): 1721-1729.

- 645 Klijn, C., S. Durinck, E. W. Stawiski, P. M. Haverty, Z. Jiang, H. Liu, J. Degenhardt, O.
646 Mayba, F. Gnad, J. Liu, G. Pau, J. Reeder, Y. Cao, K. Mukhyala, S. K. Selvaraj, M. Yu, G. J.
647 Zynda, M. J. Brauer, T. D. Wu, R. C. Gentleman, G. Manning, R. L. Yauch, R. Bourgon, D.
648 Stokoe, Z. Modrusan, R. M. Neve, F. J. de Sauvage, J. Settleman, S. Seshagiri and Z. Zhang
649 (2015). "A comprehensive transcriptional portrait of human cancer cell lines." Nat. Biotechnol.
650 **33**(3): 306-312.
- 651 Labeit, J., J. Shun and G. E. Belloch (2017). "Parallel lightweight wavelet tree, suffix array
652 and FM-index construction." J. Discrete Algorithms **43**: 2-17.
- 653 Levy-Sakin, M., S. Pastor, Y. Mostovoy, L. Li, A. K. Y. Leung, J. McCaffrey, E. Young, E.
654 T. Lam, A. R. Hastie, K. H. Y. Wong, C. Y. L. Chung, W. Ma, J. Sibert, R. Rajagopalan, N. Jin,
655 E. Y. C. Chow, C. Chu, A. Poon, C. Lin, A. Naguib, W.-P. Wang, H. Cao, T.-F. Chan, K. Y.
656 Yip, M. Xiao and P.-Y. Kwok (2019). "Genome maps across 26 human populations reveal
657 population-specific patterns of structural variation." Nat. Commun. **10**(1): 1025.
- 658 Li, Y., N. D. Roberts, J. A. Wala, O. Shapira, S. E. Schumacher, K. Kumar, E. Khurana, S.
659 Waszak, J. O. Korbel, J. E. Haber, M. Imielinski, P. S. V. W. Group, J. Weischenfeldt, R.
660 Beroukhi, P. J. Campbell and P. Consortium (2020). "Patterns of somatic structural variation in
661 human cancer genomes." Nature **578**(7793): 112-121.
- 662 MacDonald, J. R., R. Ziman, R. K. C. Yuen, L. Feuk and S. W. Scherer (2014). "The
663 Database of Genomic Variants: a curated collection of structural variation in the human
664 genome." Nucleic Acids Res. **42**(Database issue): D986-992.
- 665 Nielsen, R., J. S. Paul, A. Albrechtsen and Y. S. Song (2011). "Genotype and SNP calling
666 from next-generation sequencing data." Nature Reviews Genetics **12**(6): 443-451.
- 667 Olney, K. C., S. M. Brotman, V. Valverde-Vesling, J. Andrews and M. A. Wilson "Aligning
668 RNA-Seq reads to a sex chromosome complement informed reference genome increases ability
669 to detect sex differences in gene expression".
- 670 Paten, B., A. Novak and D. Haussler (2014). "Mapping to a Reference Genome Structure."
671 ArXiv e-prints: 1-26.
- 672 Paten, B., A. M. Novak, J. M. Eizenga and E. Garrison (2017). "Genome graphs and the
673 evolution of genome inference." Genome Res. **27**(5): 665-676.
- 674 Patro, R., S. M. Mount and C. Kingsford (2014). "Sailfish enables alignment-free isoform
675 quantification from RNA-seq reads using lightweight algorithms." Nat. Biotechnol. **32**(5): 462-
676 464.
- 677 Pedersen, B. S., P. J. Bhetariya, J. Brown, S. N. Kravitz, G. Marth, R. L. Jensen, M. P.
678 Bronner, H. R. Underhill and A. R. Quinlan (2020). "Somali: rapid relatedness estimation for
679 cancer and germline studies using efficient genome sketches." Genome Med **12**(1): 62.

- 680 Pedersen, B. S. and A. R. Quinlan (2017). "Who's Who? Detecting and Resolving Sample
681 Anomalies in Human DNA Sequencing Studies with Peddy." Am. J. Hum. Genet. **100**(3): 406-
682 413.
- 683 Poplin, R., P.-C. Chang, D. Alexander, S. Schwartz, T. Colthurst, A. Ku, D. Newburger, J.
684 Dijamco, N. Nguyen, P. T. Afshar, S. S. Gross, L. Dorfman, C. Y. McLean and M. A. DePristo
685 (2018). "A universal SNP and small-indel variant caller using deep neural networks." Nat.
686 Biotechnol. **36**(10): 983-987.
- 687 Roe, D. and R. Kuang "Accurate and Efficient KIR Gene and Haplotype Inference from
688 Genome Sequencing Reads with Novel K-mer Signatures."
- 689 Shajii, A., D. Yorukoglu, Y. William Yu and B. Berger (2016). "Fast genotyping of known
690 SNPs through approximate k-mer matching." Bioinformatics **32**(17): i538-i544.
- 691 Shen, Shen and Kidd (2020). "Rapid, Paralog-Sensitive CNV Analysis of 2457 Human
692 Genomes Using QuicK-mer2." Genes **11**(2): 141.
- 693 Shen, F. and J. M. Kidd (2020). "Rapid, Paralog-Sensitive CNV Analysis of 2457 Human
694 Genomes Using QuicK-mer2." Genes **11**(2).
- 695 Sherman, R. M., J. Forman, V. Antonescu, D. Puiu, M. Daya, N. Rafaels, M. P. Boorgula, S.
696 Chavan, C. Vergara, V. E. Ortega, A. M. Levin, C. Eng, M. Yazdanbakhsh, J. G. Wilson, J.
697 Marrugo, L. A. Lange, L. K. Williams, H. Watson, L. B. Ware, C. O. Olopade, O. Olopade, R.
698 R. Oliveira, C. Ober, D. L. Nicolae, D. A. Meyers, A. Mayorga, J. Knight-Madden, T. Hartert,
699 N. N. Hansel, M. G. Foreman, J. G. Ford, M. U. Faruque, G. M. Dunston, L. Caraballo, E. G.
700 Burchard, E. R. Bleecker, M. I. Araujo, E. F. Herrera-Paz, M. Campbell, C. Foster, M. A. Taub,
701 T. H. Beaty, I. Ruczinski, R. A. Mathias, K. C. Barnes and S. L. Salzberg (2019). "Assembly of a
702 pan-genome from deep sequencing of 910 humans of African descent." Nat. Genet. **51**(1): 30-35.
- 703 Sherman, R. M. and S. L. Salzberg (2020). "Pan-genomics in the human genome era." Nature
704 Reviews Genetics.
- 705 Sudmant, P. H., T. Rausch, E. J. Gardner, R. E. Handsaker, A. Abyzov, J. Huddleston, Y.
706 Zhang, K. Ye, G. Jun, M. H.-Y. Fritz, M. K. Konkel, A. Malhotra, A. M. Stütz, X. Shi, F. P.
707 Casale, J. Chen, F. Hormozdiari, G. Dayama, K. Chen, M. Malig, M. J. P. Chaisson, K. Walter,
708 S. Meiers, S. Kashin, E. Garrison, A. Auton, H. Y. K. Lam, X. J. Mu, C. Alkan, D. Antaki, T.
709 Bae, E. Cerveira, P. Chines, Z. Chong, L. Clarke, E. Dal, L. Ding, S. Emery, X. Fan, M. Gujral,
710 F. Kahveci, J. M. Kidd, Y. Kong, E.-W. Lameijer, S. McCarthy, P. Flicek, R. A. Gibbs, G.
711 Marth, C. E. Mason, A. Menelaou, D. M. Muzny, B. J. Nelson, A. Noor, N. F. Parrish, M.
712 Pendleton, A. Quitadamo, B. Raeder, E. E. Schadt, M. Romanovitch, A. Schlattl, R. Sebra, A. A.
713 Shabalina, A. Untergasser, J. A. Walker, M. Wang, F. Yu, C. Zhang, J. Zhang, X. Zheng-Bradley,
714 W. Zhou, T. Zichner, J. Sebat, M. A. Batzer, S. A. McCarroll, C. Genomes Project, R. E. Mills,
715 M. B. Gerstein, A. Bashir, O. Stegle, S. E. Devine, C. Lee, E. E. Eichler and J. O. Korbel (2015).
716 "An integrated map of structural variation in 2,504 human genomes." Nature **526**(7571): 75-81.
- 717 Sun, C. and P. Medvedev (2019). "Toward fast and accurate SNP genotyping from whole
718 genome sequencing data for bedside diagnostics." Bioinformatics **35**(3): 415-420.

719 Trost, B., S. Walker, Z. Wang, B. Thiruvahindrapuram, J. R. MacDonald, W. W. L. Sung, S.
720 L. Pereira, J. Whitney, A. J. S. Chan, G. Pellecchia, M. S. Reuter, S. Lok, R. K. C. Yuen, C. R.
721 Marshall, D. Merico and S. W. Scherer (2018). "A Comprehensive Workflow for Read Depth-
722 Based Identification of Copy-Number Variation from Whole-Genome Sequence Data." Am. J.
723 Hum. Genet. **102**(1): 142-155.

724 Webster, T. H., M. Couse, B. M. Grande, E. Karlins, T. N. Phung, P. A. Richmond, W.
725 Whitford and M. A. Wilson (2019). "Identifying, understanding, and correcting technical
726 artifacts on the sex chromosomes in next-generation sequencing data." Gigascience **8**(7).

727 Wood, D. E., J. Lu and B. Langmead (2019). "Improved metagenomic analysis with Kraken
728 2." Genome Biol. **20**(1): 257.

729 Xia, Y., Y. Liu, M. Deng and R. Xi (2019). "Detecting virus integration sites based on
730 multiple related sequencing data by VirTect." BMC Med. Genomics **12**(Suppl 1): 19.

731 Yang, X., W.-P. Lee, K. Ye and C. Lee (2019). "One reference genome is not enough."
732 Genome Biol. **20**(1): 104.

733