

1 **Full Title**

2 Demonstrating the utility of flexible sequence queries
3 against indexed short reads with FlexTyper
4

5 **Short Title**

6 FlexTyper: Enabling flexible queries against indexed short
7 read sequences
8
9

10 **Phillip A. Richmond^{1*}, Alice M. Kaye^{1*}, Godfrain Jacques Kounkou^{1*}, Tamar V. Av-**
11 **Shalom¹, Wyeth W Wasserman¹**
12

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

15 *These authors contributed equally.
16

17 **Corresponding Author**

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

19 Abstract

20 Across the life sciences, processing next generation sequencing data commonly relies upon
21 a computationally expensive process where reads are mapped onto a reference sequence. Prior to
22 such processing, however, there is a vast amount of information that can be ascertained from the
23 reads, potentially obviating the need for processing, or allowing optimized mapping approaches to
24 be deployed. Here, we present a method termed FlexTyper which facilitates a “reverse mapping”
25 approach in which high throughput sequence queries, in the form of kmer searches, are run against
26 indexed short-read datasets in order to extract useful information. This reverse mapping approach
27 enables the rapid counting of target sequences of interest. We demonstrate FlexTyper’s utility for
28 recovering depth of coverage, and accurate genotyping of SNP sites across the human genome.
29 We show that genotyping unmapped reads can correctly inform a sample’s population, sex, and
30 relatedness in a family setting, which can be used to inform optimized downstream analysis
31 pipelines. Detection of pathogen sequences within RNA-seq data was sensitive and accurate,
32 performing comparably to existing methods with increased flexibility. The long-term adoption of
33 the “reverse mapping” approach represented by FlexTyper will be enabled by more efficient
34 methods for FM-index generation and biology-informed collections of reference queries. In the
35 long-term, selection of population-specific references or weighting of edges in pan-population
36 reference genome graphs will be enabled by the FlexTyper reverse mapping approach. FlexTyper
37 is available at <https://github.com/wassermanlab/OpenFlexTyper>.

38

39 Author Summary

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

61

62 Introduction

63 Short-read DNA sequencing enables diverse molecular investigations across life science
64 applications spanning from medicine to agriculture. Obtaining useful information from sequencing
65 datasets typically involves either performing *de novo* assembly, or mapping the data against one
66 or more reference genomes. The process of mapping sequencing reads (short pieces of DNA read-
67 outs from the DNA sequencer) against reference genomes, or a collection of reference genomes,
68 is made computationally tractable by indexing the reference sequences, commonly performed with
69 a Burrows Wheeler transform or FM-index. Several data analysis pipelines, whether they focus on
70 quantification (e.g. observed gene expression in RNA sequencing data), or identifying sequence
71 differences between a sample and a reference genome (e.g. genotyping), leverage reference
72 genome mapping as a primary analysis component.

73 While the status quo has been to utilize linear representations of reference genomes, a
74 transition away from a single haploid reference genome is inevitable (Yang et al. 2019; Ballouz,
75 Dobin, and Gillis 2019). This transition is supported by several factors. A large amount of
76 structural variation exists between human populations (Feuk, Carson, and Scherer 2006;
77 MacDonald et al. 2014; Levy-Sakin et al. 2019). A recent study focusing on ~1000 individuals of
78 African descent identified nearly 200 million bases missing from the most recent reference genome
79 (Sherman et al. 2019). Static linear reference genomes which do not capture these large differences
80 between populations impose challenges for accurate genotyping (Ballouz, Dobin, and Gillis 2019;
81 Yang et al. 2019), with implications in medicine and association studies. An alternative to choosing
82 from a collection of population-specific reference genomes is to use emerging graph genome
83 approaches to unite the data (Dilthey et al. 2015). As highlighted in a review by (Paten et al. 2017),
84 in either approach, a key challenge in the future will be to determine the most appropriate reference
85 genome(s), or path(s) through a graph genome, to maximize genotyping performance. Knowledge
86 of distributed single nucleotide polymorphisms (SNPs) genotypes across the genome can be used
87 to guide such choices.

88 Currently, the primary approach for identifying SNP genotypes across the genome utilizes
89 computationally expensive reference-based read mapping and variant calling strategies (Nielsen
90 et al. 2011). Inferring ancestry from specific, population-discriminating SNPs can be performed
91 rapidly with the recently published tool Peddy, which uses fewer than 25,000 SNPs to identify

92 ancestry through principal component analysis (Pedersen and Quinlan 2017). Previous work
93 showed that it is possible to genotype predefined SNPs from unmapped sequence data,
94 circumventing the read mapping and variant calling process (Dolle et al. 2017; Sun and Medvedev
95 2019; Shajii et al. 2016). Some approaches focus on kmer (short sequences of length k) hashing
96 and matching to predefined target kmers to perform genotyping of known SNPs, as demonstrated
97 in the VarGeno and LAVA frameworks (Sun and Medvedev 2019; Shajii et al. 2016). These
98 approaches are fast, but rely upon indexes of kmers extracted from the reference genome and SNP
99 databases, thus reducing their flexibility for kmers of different length and source. A separate
100 approach is taken by Dolle et al., wherein the entire 1000 Genomes dataset is compressed into an
101 FM-index and queried with kmers spanning polymorphic sites, thus demonstrating the utility of
102 scanning unmapped reads for predefined kmers of interest. The “reverse mapping” highlighted in
103 their approach was applied to aggregated data, but the concept can be extended to the analysis of
104 individual genomes if implemented in a flexible way for diverse types of queries.

105 Within the paradigm of indexing reads and performing reverse mapping, other useful
106 operations can be performed with increased utility, especially in cases with a diverse set of
107 informative sequences. One example of this is within RNA sequencing (RNA-seq), where analysis
108 of cancer RNA-seq datasets can reveal the presence of viral pathogens within patient data (Klijn
109 et al. 2015). Several tools have been developed to specifically detect these viral pathogens from
110 sequencing data including viGEN (Bhuvaneshwar et al. 2018) and VirTect (Xia et al. 2019).
111 However, they are hampered by a computationally expensive iterative mapping procedure which
112 first maps against the human reference genome and then subsequently maps against viral genome
113 collections. Other methods, such as Centrifuge (Kim et al. 2016) and Kraken2 (Wood, Lu, and
114 Langmead 2019), rely upon kmer searches against large viral and bacterial databases. Both of these
115 methods are powerful, but come with drawbacks of flexibility and reliance upon phylogenetic
116 relationships between target sequences. Specifically, they require re-indexing of search databases
117 for different query lengths or when the target sequences change. Nevertheless, these tools are
118 broadly used and thus serve as good comparators for efficacy, as they have both been demonstrated
119 to have utility in detecting viral pathogens within cancer RNA-seq datasets by examining kmer
120 content. (<https://www.sevenbridges.com/centrifuge/>).

121 Combining the current drive to decrease our reliance upon linear reference genomes, and
122 the wealth of demonstrated utility of reverse mapping approaches, we developed FlexTyper.

123 FlexTyper is a computational framework which enables the flexible indexing and querying of raw
124 next generation sequencing reads. We show example usage scenarios for FlexTyper by
125 demonstrating the high accuracy of reference-free genotyping of SNPs in single samples, and the
126 ability to identify foreign pathogen sequences within short-read datasets. We hope the flexibility
127 afforded by the framework underpinning FlexTyper will fuel the emerging trend away from the
128 necessity for a static reference genome that currently lay at the heart of the majority of genomic
129 analysis tools.

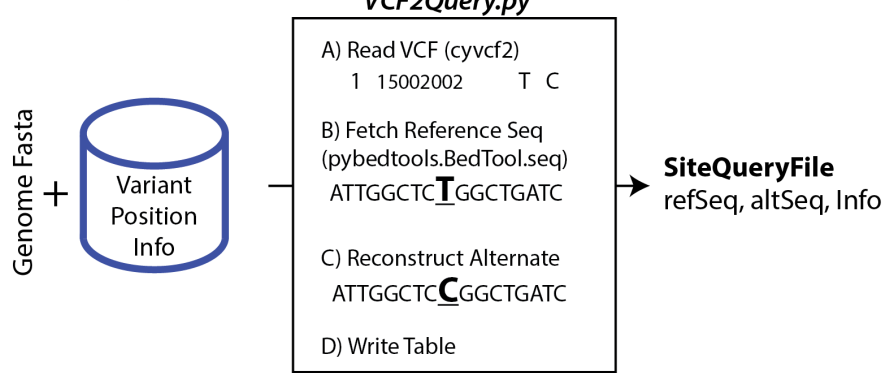
130 Design and Implementation

131 Overview of FlexTyper

132 Usage can be broken down into three steps: 1) query generation, 2) indexing the raw reads,
133 and 3) querying against the FM-index (Figure 1). For query generation, we allow for both custom
134 user query generation, as well as pre-constructed queries from useful databases, such as
135 CytoScanHD array probe queries. Custom queries designed to capture genomic loci can be
136 generated by pairing a user-provided VCF (format v4.3) with a reference genome fasta file. For
137 the capture of potential pathogen sequences, we also allow query generation from one or more
138 fasta files. The files produced from query generation are used as input for subsequent index query
139 operations. The second step is the production of an FM-index from a set of short-read sequences
140 in fastq format. This process includes reverse-complementing the entire read file, and
141 concatenating the transformed reads with the original set. This is done in order to prevent the need
142 to scan for the reverse complement of the query kmers . The third step is the core FlexTyper search
143 algorithm which takes the query input file, generates search kmers, and scans the FM-index for
144 matches. This step creates an output with matching format to the input file, with appended counts
145 of matching reads for each query. A detailed breakdown of these three components is described
146 below.

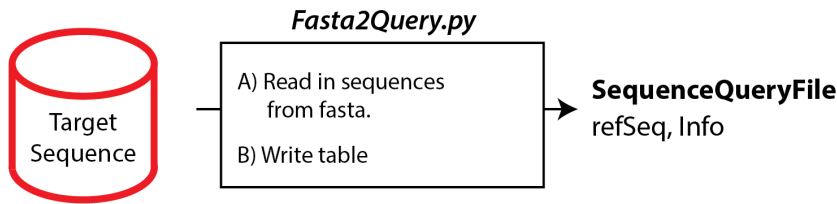
147

1) Query Generation



Directed Usage

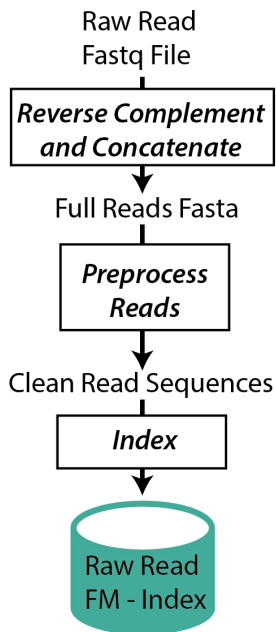
Querying for specific genomic sites and recovering counts of reads over those sites. E.g. variant detection given target sites/VCF.



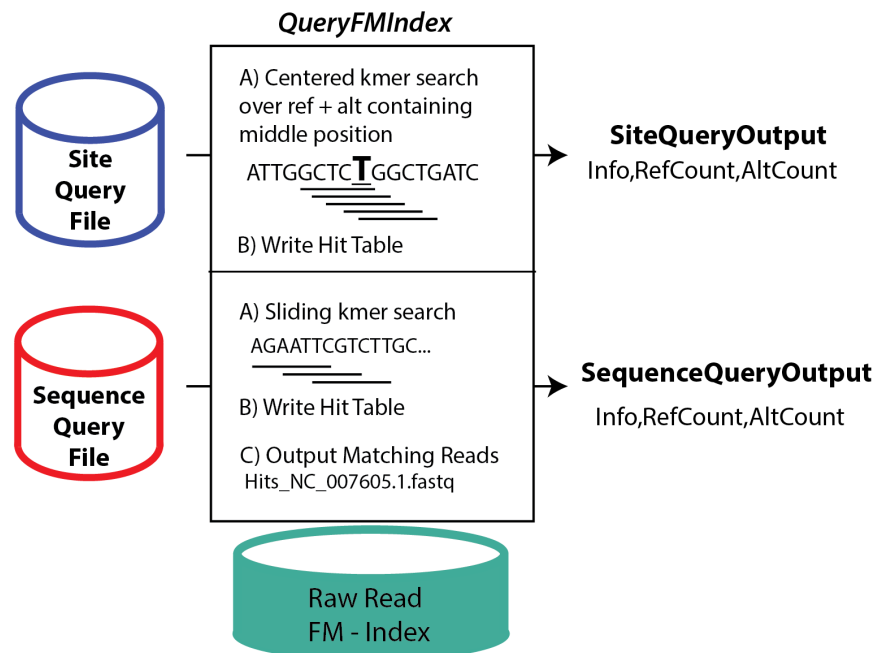
Directed Usage

Querying specific sequences, even if the length of the query is larger than the length of the search parameter 'k'. E.g. pathogen detection.

2) FM Indexing Reads



3) Query Against FM-Index



148

149 Fig 1 - Overview of FlexTyper

150 FlexTyper has three primary components: query generation, FM indexing reads, and querying
 151 against the FM-index. Query generation includes the capacity to translate VCF files into query
 152 files given a reference genome file (Genome Fasta), or to directly create queries from fasta
 153 sequences including pathogen genome sequences. Modules VCF2Query.py and Fasta2Query.py
 154 facilitate this process. The second component involves creating an FM-index of the raw reads,
 155 after reverse complementing and concatenating the read set and performing optional
 156 preprocessing steps. The third component executes the queries against the FM-index to produce
 157 output files with counts of reference and alternate sequences within the query files.

158

159 Query generation

160 FlexTyper supports flexible query generation giving users the capacity to query for any
161 target sequence or allele within their read dataset. Query files can be generated from an input fasta
162 and VCF file (VCF2Query.py), or directly from a fasta file (Fasta2Query.py). Potentially useful
163 queries, including those presented here, are provided online and include all sites from the
164 CytoScanHD chromosomal microarray, and ancestry discriminating sites (Pedersen and Quinlan
165 2017). These predefined query sets are available through git-lfs in the online FlexTyper github
166 repository (<https://github.com/wassermanlab/OpenFlexTyper>). If users wish to directly query a
167 short-read dataset with a set of predetermined kmers, they can provide the kmers as a fasta file and
168 set the k parameter to the length of the kmers in the file.

169 FM-index creation

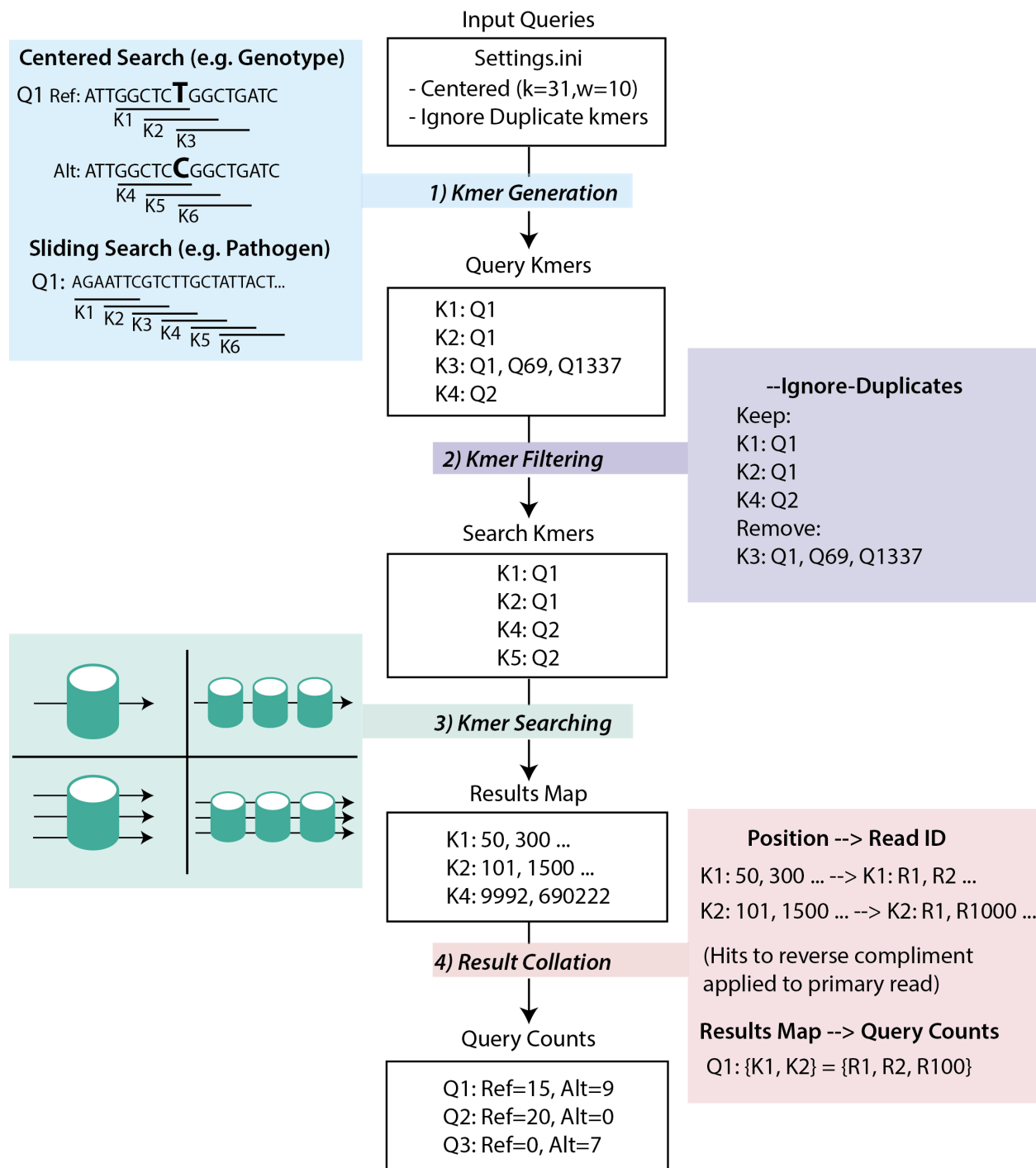
170 Generating the FM-index from short-read sequencing datafiles is performed in two steps;
171 preprocessing and indexing. The focus of our work is not on the algorithms used to construct the
172 FM-index, and hence we use two existing utilities to generate a compatible FM-index for
173 FlexTyper. The toolkit Seqtk is used for reformatting compressed fastq files by removing quality
174 scores and non-sequence information to create a sequence-only fasta format, and append this with
175 the reverse complement of the reads. The output fasta file is then processed using the SDSL-Lite
176 library to generate the FM-index. SDSL builds a suffix array that is used to generate the BWT of
177 the input string, which is then compressed using a wavelet tree and subsampled. The resulting
178 compressed suffix array is streamed to a binary index file. As the memory requirements for
179 indexing large files can be burdensome, we support an option to split the input file and index each
180 chunk of reads independently. Downstream search operations support the use of multiple indexes.
181

182 Query against FM-index

183 Querying the FM-index for user selected sequences can be conceptually divided into four
184 steps: 1) kmer generation; 2) kmer filtering; 3) kmer searching; and 4) result collation (Figure 2).

185 There are two primary methods of kmer generation for a query; a centered search where the middle
186 position of the query is included in all kmers, and a sliding search which starts at one end of the
187 query and uses a sliding window approach to generate the kmers (Figure 2). Centered search can
188 be used for genotyping or estimating coverage over a single position, and the sliding search can be
189 used to count reads which match to any part of a query sequence. The *--ignore-duplicates*
190 parameter filters query kmers by ignoring kmers that occur in multiple query sequences. After
191 filtration, the kmers are searched for within the FM-index using C++ multithreading and
192 asynchronous programming, using either a single thread on a single index, multiple threads on a
193 single index, a single thread on multiple indexes, or multiple threads on multiple indexes (Figure
194 2). Importantly, asynchronous programming allows the number of threads used during searching
195 to be increased beyond the number of available CPUs. The output from this search process is a
196 collated results map containing the positions of each kmer within the FM-index. These positions
197 are translated to read IDs, and finally collapsed into query counts using the kmer-to-query
198 mapping. Importantly, if multiple kmers from the same query hit the same read, they are recorded
199 as a single count at the query level. For cases of multiple indexes being searched in parallel, the
200 kmer searching and assignment to the query count is performed independently and then merged to
201 produce a final query count table.

202



203

204 Figure 2 - Query Search Workflow

205 Workflow for query search against the FM-index, starting with input queries and settings defined
 206 in Settings.ini file. In this example, it sets a centered search with ignoring duplicate kmers enabled.

207 1) Kmer generation has two modes, centered search and sliding search. For a centered search, the
 208 position of interest lies in the middle of the query, and kmers are designed to overlap that central
 209 position with defined length (k) and step (w). 2) If the ignore-duplicates option is set, kmers

210 collated from the query set are filtered to remove any kmers which were found in multiple query
211 sequences. 3) The filtered kmers are then searched for within the FM-index (left two panels) or
212 multiple indexes (right two panels) of the read set. This can be done using single (top two panels)
213 or multiple (bottom two panels) threads. 4) The results corresponding to a position within the FM-
214 index are then translated back into reads, with hits on reverse complement reads assigned to the
215 primary read, and collapsed into a set for each query. The final counts are reported per query.

216

217

218

219 Post-processing of results into downstream formats

220 The output tables from the search process for genotyping can be translated into useful formats
221 for downstream analysis using the fmformatter scripts

222 (<https://github.com/wassermanlab/OpenFlexTyper/tree/master/fmformatter>). Currently, there is
223 the capacity to output genotype calls in VCF, 23andMe, or Ancestry.com format. Genotype calls
224 are derived here using a basic approach which assigns genotypes given a minimum read count
225 parameter as follows:

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

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

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

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

231

232

233

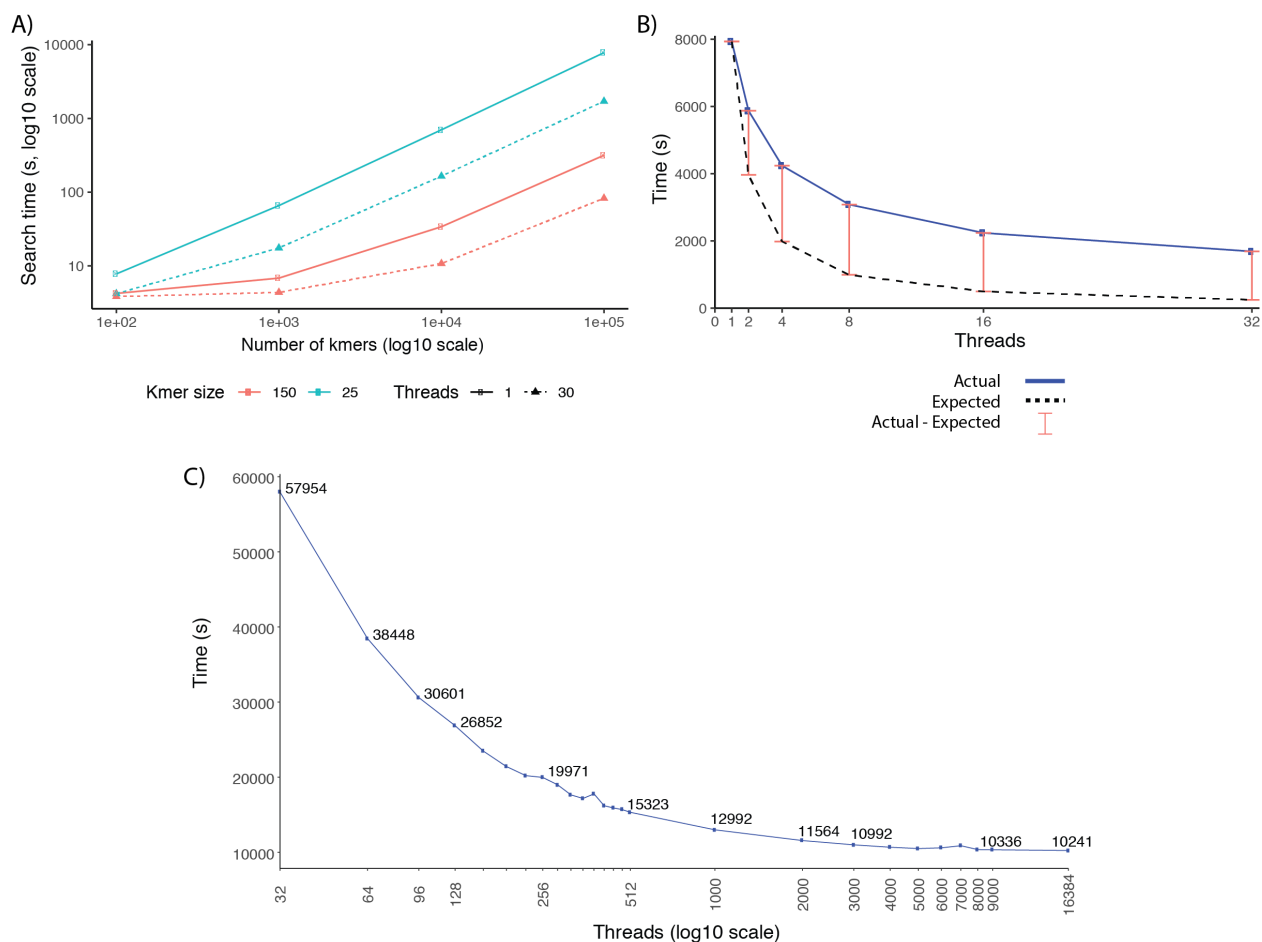
234

235 Results

236 Performance metrics for indexing and querying

237 We used a human whole genome sequencing (WGS) sample to demonstrate the indexing
238 and querying capacities of FlexTyper. Our indexing strategy utilizes open source tools to build the
239 FM-index on a high memory CPU, with at least 1000GB of RAM. While index creation
240 optimization was not the focus of this work, indexing is feasible on standard systems with ~256GB
241 of RAM, as long as the input read dataset is smaller than ~20GB (Supplemental Table S1). Since
242 our search function allows for multiple separate indexes, we incorporated the ability to sub-divide
243 larger read sets into multiple smaller read sets that can be indexed in parallel. For querying, we
244 generated a set of kmers designed from probes on the CytoScanHD Illumina genotyping
245 microarray with a centered search process (Figure 2) for varying kmer lengths (Figure 3A). The
246 CytoScanHD genotyping microarray, chosen for its broad usage in the field of human genetics,
247 has probe sequences designed to uniquely detect well-characterized SNPs. There is a noticeable
248 benefit from multithreading FlexTyper, which we demonstrated by isolating the kmer searching
249 process across 1 to 32 threads (Figure 3A). As the number of threads increases, we observe a
250 continuous decrease in search time, and by comparing between observed and expected
251 performance, the performance advantage gained from additional threads does not plateau at 32
252 threads (Figure 3B). As the software is written using asynchronous programming, we tested the
253 upper bound on allocated data threads given a fixed set of 32 CPUs on a single machine with
254 256GB of RAM. For this analysis, we used an extended set of queries from the CytoScanHD SNP
255 set, for a total of ~6.4 million kmers. We increased the threads from 32 to 512 stepping by 32 and
256 while we do see a decrease on the improvement in speed, there is still a benefit of additional threads
257 (Figure 3C). To see where the benefit of increased threads plateaus, we increased threads from
258 1000 to 16,384 and witnessed little speed increase (<5 minutes) between 5000 and 16,384 threads
259 (Figure 3C). Thus, we define the upper bound on data threads for a machine with 32 CPUs and
260 250GB of RAM to be ~5000 data threads, for a query set of ~6.4 million kmers. It is possible that
261 higher thread counts may improve performance for larger query sets and more powerful
262 computers. Lastly, to highlight the clear advantage over non-indexed methods, we compared

263 FlexTyper to popular non-indexed algorithms achieving a decrease in search time by roughly three
 264 orders of magnitude when using FlexTyper (Supplemental Table S2).
 265



266
 267 **Figure 3 - Search speeds for FlexTyper**
 268 A) FlexTyper search time with kmers of size 25 (blue) and 150 (red), increasing in number from
 269 10-100,000, using one (solid) or 30 (dashed) threads. B) Increasing the number of threads from 1
 270 to 32, for 100,000 kmers of length 25 (solid blue line). Expected values calculated by dividing
 271 single thread time by the additional number of threads (dashed black line), with difference between
 272 actual and expected plotted (red vertical bar). C) Hyperthreading results for the time (in seconds)
 273 vs. thread counts from 32 up to 16,384 (log₁₀ scaled x-axis)

274 **Genomic coverage and genotype detection within human WGS data**

275 Knowing whether a given kmer is present or absent from a human WGS datafile (in this
 276 instance genome Illumina short-read, paired-end data) can have utility for estimating the depth of
 277 coverage for a target region and genotyping SNPs. FlexTyper has the capacity to compute depth
 278 of coverage or genotype SNPs from WGS data for both predefined and user-supplied loci. We

279 demonstrate this capacity for genomic sites using the probe sequences from the CytoScanHD
280 microarray, as well as a subset of previously collated population discriminating SNPs (Pedersen
281 and Quinlan 2017). Using these loci, we created query files with a reference and alternate query
282 sequence centered on the biallelic site (Supplemental Methods).

283 We first sought to test the read recovery capacity of FlexTyper compared to an alignment
284 based method which we call BamCoverage. The BamCoverage method involves mapping the
285 reads to the reference genome, and then extracting per-base read coverage over a specific reference
286 coordinate. BamCoverage utilizes the pysam package to extract read pileup over positions defined
287 by the FlexTyper input query file (Supplemental methods). Using the CytoScanHD SNP set, we
288 found a high concordance between the read counts from FlexTyper and the depth of coverage from
289 aligned reads (Figure 4A). The vast majority, 780,178/797,653 or 97.8%, of sites differed by less
290 than 10 between FlexTyper and BamCoverage (Figure 4B). This discrepancy is similar for both
291 reference and alternate alleles, which is important since most genotyping models assume relative
292 contributions of observed alleles for genotype calling. There were 16,282 sites with a delta, ($\Delta =$
293 FlexTyper - BamCoverage), greater than 10, and 4,256 sites with a delta greater than 100. We
294 manually investigated a few of these sites which were overcounted by FlexTyper by more than
295 100 and found that they are being overcounted due to kmers mapping to multiple possible
296 locations. Comparing these over-counted hits with delta greater than 100 to previously defined
297 repeat regions shows that 4189/4256 or 98.4% of the overcounted sites overlapped with predefined
298 repeats (Trost et al. 2018). The uniqueness of kmers is important for accurate read counting, thus
299 it is recommended to filter such regions when using FlexTyper for genotyping or depth profiling.
300 Lastly, by examining the recovery of reads across the chromosome between FlexTyper and the
301 read alignment approach, it's clear that FlexTyper can accurately capture relative sequence
302 abundance with relevance to copy number variant calling applications (Fig 4D).

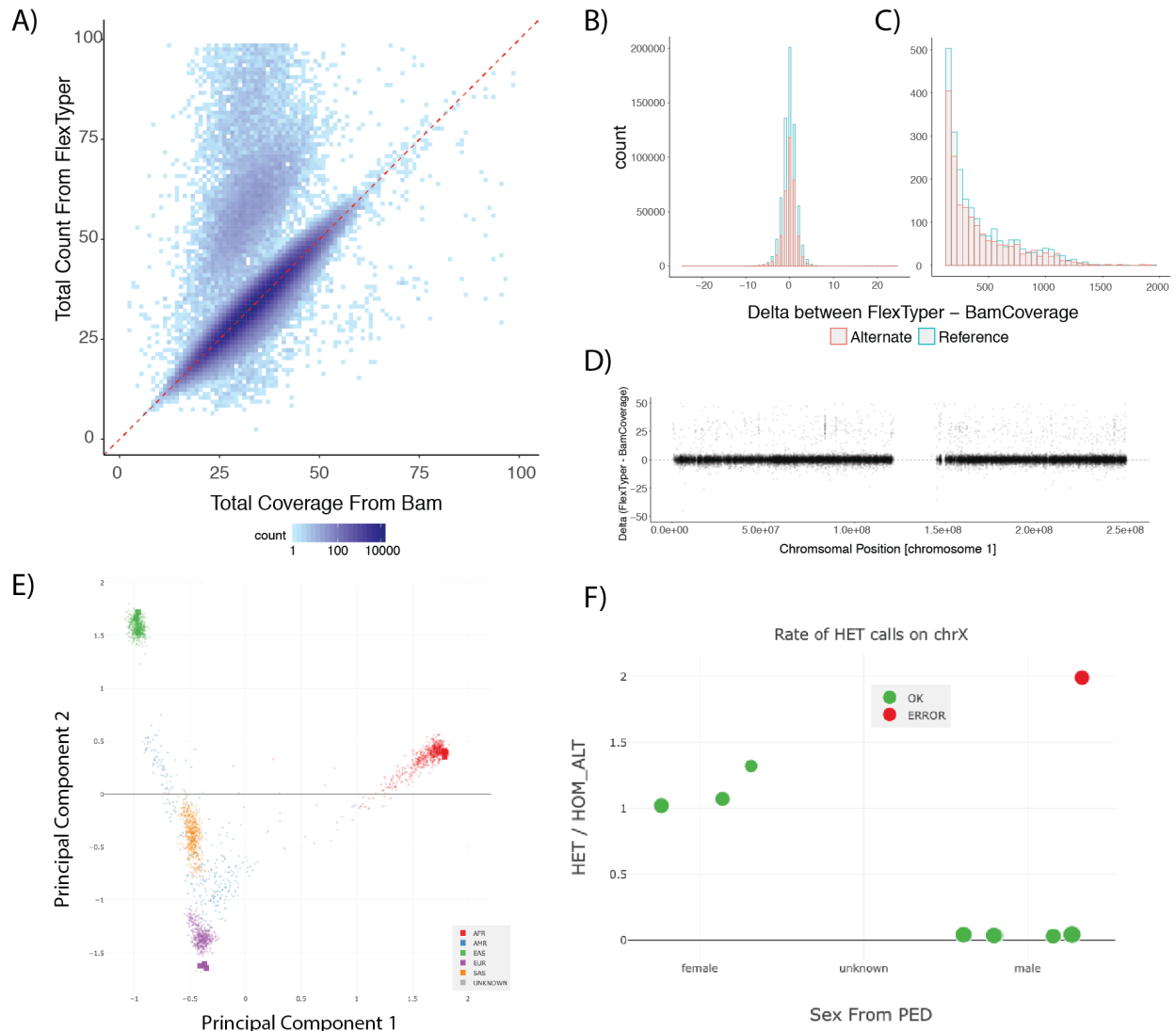
303 Next, we investigated whether FlexTyper can accurately recover genotypes at the SNP sites
304 profiled from the chromosomal microarray. The genotyping approach we use leverages a minimum
305 count from the reference and alternate allele to assign heterozygous, homozygous alternate, or
306 homozygous reference genotypes (Supplemental Methods). We applied this basic genotyping
307 algorithm to both FlexTyper and BamCoverage counts to produce a VCF file. These genotypes
308 were compared to an alternate pipeline which uses reference-based mapping and sophisticated
309 variant calling using DeepVariant (Poplin et al. 2018). For the 797,653 SNPs on the CytoScanHD

310 microarray, all three methods agree on 99.2% (791,063/797,653) of the sites. For the sites where
311 there was disagreement, we see an overlap with the repeat regions of 5004/5586 or 89.6%,
312 affirming that these repeat regions are responsible for the majority of discordant genotypes. We
313 further demonstrate the accuracy of these genotypes by indexing nine WGS samples from the
314 Polaris project representing diverse populations including three African, three Southeast Asian,
315 and three European individuals (S. Chen et al., n.d.). After indexing, we queried the samples for
316 population discriminating sites and then genotyped the output table to produce a VCF file. The
317 output VCFs were then used within the Peddy tool, and a principal component analysis was
318 performed to predict the ancestry of the samples (Pedersen and Quinlan 2017). In all nine cases
319 the population was correctly determined, as well as the relatedness inference for the three trios
320 (Figure 4E, Figure S1). Interestingly, we observed a discrepancy between the listed sex for the
321 child of the European trio, individual HG01683, and the inferred sex from FlexTyper and Peddy
322 (Figure 4F). We followed up on this observation and revealed that the individual is not an XY
323 male, but rather an XXY individual. Taken together, FlexTyper has the capacity to provide
324 accurate counts of observed reads matching a query sequence, with relevant utilities such as copy
325 number estimation, sample identification, ancestry typing, and sex identification.

326

327

328



329

330 Figure 4 - WGS Genotyping using FlexTyper

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

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

344 To demonstrate the capacity of FlexTyper to detect pathogens from RNA-sequencing data,
345 we generated synthetic reads from four relevant viral genomes including Epstein-Barr virus
346 (EBV), Human Immunodeficiency virus type 1 (HIV-1), and two Human Papilloma virus strains
347 68b (HPV FR751039) and 70 (HPV U21941) (Supplemental Methods). We first examined the
348 impact of various FlexTyper parameters on the recovery rate of pure, simulated read sets for each
349 of the four viruses and one human blood RNA-seq dataset from the Genome England project
350 (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-6523/samples/>) (Table S3).
351 Importantly, varying the parameters k (length of search substring) and w (step-size) change the
352 specificity and sensitivity of read recovery. When k is set to 15 (a short kmer), there are roughly 1
353 million off-target hits to the viral genomes for the pure human RNA-seq file (Table S3). Next, we
354 demonstrated that the kmer uniqueness setting only guarantees that identical kmers cannot appear
355 across queries. Thus, if query specificity is a priority then setting the w parameter to 1 will produce
356 results with the least amount of cross-query assignment. By exploring these parameters, we show
357 that all simulated reads can be recovered with parameters of 30 and 5 for k and w respectively,
358 with low off-target assignment.

359 Next, to simulate patients infected by one of the four viruses, we spiked-in simulated
360 pathogen reads with the human RNA-seq dataset. Using the optimized parameters derived above,
361 we are able to detect each virus in the patient sample even at low concentrations (Table S4). We
362 further demonstrate the capacity for FlexTyper to discriminate between spiked-in virus samples
363 by mixing the viruses at differing concentrations (read counts) within the human RNA-seq dataset
364 (Table 1). FlexTyper was run with two settings by varying the k and w parameters for increased
365 sensitivity ($k=31$, $w=5$), or increased speed ($k=100$, $w=25$). We compared these results with
366 Centrifuge, a software tool that works with unmapped short-read sequencing data by performing
367 read-length ($k=150$) kmer searches against a database of viral and bacterial genomes (Kim et al.
368 2016). For each of the datasets, FlexTyper is able to detect the contaminating pathogen sequences,
369 even in a sample (Patient_5) where we only spiked in the equivalent of a 1x coverage of the viral
370 genome, which equates to roughly 50-1150 reads depending on the size of the viral genome (Table
371 1). For comparison, Centrifuge results were manually combined for viruses of similar naming
372 schema and presented as the sum of non-unique read hits (Table 1). For each of the samples, both
373 Centrifuge and FlexTyper are capable of detecting the spiked-in pathogens. In all simulations,

374 consistent with its use of shorter kmers, FlexTyper is more sensitive in its detection capacity,
375 recovering more reads than Centrifuge. In our collation of viral hits for the Centrifuge data, we
376 observed the limitation that for HPV strains, Centrifuge utilizes a comprehensive genome database
377 with hundreds of distinct strains. Thus, retrieving a combined count of HPV sequences within a
378 sequencing dataset is nontrivial and requires collation over hundreds of viral genome hits. In
379 contrast, FlexTyper is able to detect all of the spiked in reads for these viral genomes of interest.
380 This is due to the increased flexibility of FlexTyper, which enables the user to define the relevant
381 pathogens to search for without the need for reconstructing a complex bacterial or viral database,
382 as is the case for Centrifuge. In summary, FlexTyper is more sensitive in its detection capacity
383 than Centrifuge, and the flexibility to *ad hoc* define the pathogen search space could be beneficial
384 in some applications, such as instances when the virus is a novel strain.

385

386

387

| | EBV | | | | HIV-1 | | | | | |
|-----------|-----------|---------|-----------|--------|--------|---------|-----------|----------|---------|-----------|
| Sample | Expect | FT:30,5 | FT:100,25 | C | Expect | FT:30,5 | FT:100,25 | C | | |
| Patient_1 | 1145 | 1146 | 861 | 538 | 610 | 610 | 462 | 284 | | |
| Patient_2 | 11450 | 11454 | 8290 | 5427 | 6100 | 6099 | 4359 | 2797 | | |
| Patient_3 | 114500 | 114502 | 83504 | 54352 | 61000 | 61000 | 43648 | 28206 | | |
| Patient_4 | 1145000 | 1144990 | 833955 | 543924 | 62 | 62 | 45 | 29 | | |
| Patient_5 | 1146 | 1146 | 861 | 538 | 62 | 62 | 45 | 29 | | |
| | | | | | | | | | | |
| | Total HPV | | | | U21941 | | | FR751039 | | |
| Sample | Expect | FT:30,5 | FT:100,25 | C | Expect | FT:30,5 | FT:100,25 | Expect | FT:30,5 | FT:100,25 |
| Patient_1 | 57200 | 60812 | 40930 | 2650 | 5200 | 8475 | 3755 | 52000 | 52337 | 37175 |
| Patient_2 | 52052 | 55371 | 37443 | 1470 | 52000 | 52005 | 37406 | 52 | 3366 | 37 |
| Patient_3 | 572 | 615 | 422 | 25 | 52 | 92 | 38 | 520 | 523 | 384 |
| Patient_4 | 5720 | 6084 | 4152 | 273 | 520 | 851 | 388 | 5200 | 5233 | 3764 |
| Patient_5 | 104 | 112 | 75 | 3 | 5 | 57 | 38 | 52 | 55 | 37 |

388

389 Table 1 - Performance comparison for simulated spike-in pathogens.

390 Each of the samples, (Patient_1 - Patient_5), with expected (simulated known counts) vs. observed
 391 counts for Centrifuge (C) and FlexTyper with $k=30/w=5$ (FT:30,5) and $k=100/w=25$ (FT:100,25).
 392 Each quantified viral strain includes Epstein Barr Virus (EBV), Human Immunodeficiency Virus-
 393 1 (HIV-1), total Human Papillomavirus (HPV), and two strains of HPV (type 70 and 68b). The
 394 maxOcc parameter was set to limit the number of hits to one million and non-unique kmers were
 395 allowed. For centrifuge, sub-strain HPV counts were not feasible so counts were aggregated over
 396 all papilloma viral strains in the output report file per patient.

397

398

399

400

401

402

403 Discussion

404 Here we presented FlexTyper, a flexible tool which enables exploratory analysis of short
405 read datasets without the need for alignment to a reference genome. Our framework allows for the
406 custom generation of queries, giving the user total control to perform searches relevant to the
407 problem at hand. We demonstrated three applications, including depth of coverage analysis,
408 accurate SNP genotyping, and sensitive detection of pathogen sequences. FlexTyper is available
409 for the creative use of genomics researchers.

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

419 The genotyping case study highlights how pre-alignment analysis of genome sequence data
420 can provide rapid insights into the properties of a sample. SNP genotyping was accurate across the
421 genome, allowing rapid identification of sample ancestry, sample relatedness in the trio setting,
422 and sample sex typing using Peddy (Pedersen and Quinlan 2017). Interestingly, applying Peddy to
423 the output of FlexTyper for open source trio data from the Polaris project revealed a mislabeling of
424 the sex for individual HG01683, which was reported and subsequently amended in the online
425 data repository (<https://github.com/Illumina/Polaris/wiki/HiSeqX-Kids-Cohort>). Since ancestry
426 and sex information can inform choices in downstream data processing, identifying these
427 discrepancies between labeled sex and inferred sex in a data-driven manner is a critical step of pre-
428 alignment informatics. For instance, mapping against the sample-matched sex chromosomes has
429 been shown to improve performance (Webster et al. 2019; Olney et al., n.d.). As such, using
430 FlexTyper, in combination with Peddy, on diverse datasets prior to reference-guided read
431 alignment will lead to improved results from mapping-based pipelines.

432 The importance of pathogen identification is increasingly recognized. In both cancer
433 profiling (Klijn et al. 2015) and public health studies (Gardy, Loman, and Rambaut 2015), rapid
434 determination of the presence of pathogen sequences could obviate the need for full reference
435 mapping. Some existing tools designed for viral detection in sequencing data rely upon pre-
436 indexed databases of viral and bacterial sequences, sometimes including a phylogenetic
437 relationship between genomes within the index (Xia et al. 2019; Wood, Lu, and Langmead 2019;
438 Kim et al. 2016). One such approach, Centrifuge, has been applied to cancer genomes to confirm
439 the presence of viral pathogens. We demonstrated that our approach compares favorably to
440 Centrifuge, with a more sensitive detection level, due to the ability to search for kmers shorter than
441 the read length and the advantage of fine-tuned control over the searchable database. Here we only
442 searched for viral pathogens of interest, although other specific pathogen queries could be
443 performed, such as the presence of antibiotic resistance genes within a patient RNA-seq sample.

444 We anticipate that the research community will identify diverse and creative uses for
445 “reverse mapping” analysis with FlexTyper, but a few approaches are apparent to us. It is feasible
446 to genotype complex structural variants by searching for sequences overlapping breakpoints, such
447 as those observed in a subpopulation, or events recurrently found in cancer (Li et al. 2020; Sudmant
448 et al. 2015). Within RNA-seq data, querying for exon-exon splice junctions in a rapid manner can
449 allow isoform quantification, as has been previously demonstrated (Patro, Mount, and Kingsford
450 2014; Bray et al. 2016). Further, a recent report showed the utility of kmer-counting methods in
451 resolving copy number variants within paralogous loci and genes (Shen, Shen, and Kidd 2020).
452 Another group showed the advantage of examining depth of coverage at specific sites across the
453 paralogous genes in Spinal Muscular Atrophy (X. Chen et al. 2020) As FlexTyper is well suited
454 for specific sequence recovery operations, scanning with preselected query sequences such as
455 defined by these studies can enable rapid detection (X. Chen et al. 2020). All of these proposed
456 applications help tackle challenges which are currently a burden for traditional reference-based
457 mapping approaches.

458 We focused this report on the utility of kmer searches against indexed read sets, but
459 recognize that speed and computational resources are an important consideration for adoption of
460 the method. One obvious (but transient) constraint on the utility of FlexTyper is the generation of
461 the FM-index for the sequencing reads. As the FM-index is critical to many aspects of genome-
462 scale sequence analyses, there are diverse efforts to develop novel indexing strategies, such as

463 optimizing FM-index construction using GPUs (Chacón et al. 2015) and creating efficient
464 construction algorithms (N. Chen, Li, and Lu 2018; Labeit, Shun, and Blelloch 2017). Further, the
465 nature of the mapping procedure holds promise with massive parallelization approaches, including
466 those involving GPU acceleration (Hung et al. 2018). Moving forward, accelerations to the FM-
467 index generation and reverse mapping approach will result in faster genomic analysis pipelines
468 than is currently possible with alignment based methods.

469 Looking to the future, we see the kmer-searching approach of FlexTyper as having great
470 utility when used in conjunction with emergent graph-based representations of the reference
471 genome (Kehr et al. 2014; Paten et al. 2017; Kaye 2016). Whether users seek to select a population
472 specific reference graph as the basis for read mapping, or to introduce Bayesian priors (edge
473 weighting) within a pan-population reference graph, knowledge of population markers spanning
474 chromosomes will be required to inform the processes. Furthermore, it is our expectation that
475 graph-based mapping methods will ultimately use read-based FM-indices, as indexing the
476 reference graph imposes restrictions on the graph structures that can be used and the types of
477 variations that can be incorporated (Ghaffaari and Marschall 2019; Paten, Novak, and Haussler
478 2014). As the graph-based algorithms mature, approaches such as FlexTyper which enable reverse
479 mapping of sequences against a set of indexed reads will be instrumental in the initial steps of
480 genome analysis pipelines, and in the resolution of challenging regions of the genome.

481

482

483 References

- 484 Ballouz, Sara, Alexander Dobin, and Jesse A. Gillis. 2019. “Is It Time to Change the Reference
485 Genome?” *Genome Biology* 20 (1): 159.
- 486 Bhuvaneshwar, Krithika, Lei Song, Subha Madhavan, and Yuriy Gusev. 2018. “viGEN: An Open Source
487 Pipeline for the Detection and Quantification of Viral RNA in Human Tumors.” *Frontiers in*
488 *Microbiology* 9 (June): 1172.
- 489 Bray, Nicolas L., Harold Pimentel, Páll Melsted, and Lior Pachter. 2016. “Near-Optimal Probabilistic
490 RNA-Seq Quantification.” *Nature Biotechnology* 34 (5): 525–27.
- 491 Chacón, Alejandro, Santiago Marco-Sola, Antonio Espinosa, Paolo Ribeca, and Juan Carlos Moure. 2015.
492 “Boosting the FM-Index on the GPU: Effective Techniques to Mitigate Random Memory Access.”
493 *IEEE/ACM Transactions on Computational Biology and Bioinformatics / IEEE, ACM* 12 (5): 1048–
494 59.
- 495 Chen, N., Y. Li, and Y. Lu. 2018. “A Memory-Efficient FM-Index Constructor for Next-Generation
496 Sequencing Applications on FPGAs.” In *2018 IEEE International Symposium on Circuits and*

- 497 *Systems (ISCAS)*, 1–4.
- 498 Chen, Sai, Peter Krusche, Egor Dolzhenko, Rachel M. Sherman, Roman Petrovski, Felix Schlesinger,
499 Melanie Kirsche, et al. n.d. “Paragraph: A Graph-Based Structural Variant Genotyper for Short-Read
500 Sequence Data.” <https://doi.org/10.1101/635011>.
- 501 Chen, Xiao, Alba Sanchis-Juan, Courtney E. French, Andrew J. Connell, Isabelle Delon, Zoya Kingsbury,
502 Aditi Chawla, et al. 2020. “Spinal Muscular Atrophy Diagnosis and Carrier Screening from Genome
503 Sequencing Data.” *Genetics in Medicine: Official Journal of the American College of Medical
504 Genetics*, February. <https://doi.org/10.1038/s41436-020-0754-0>.
- 505 Dilthey, Alexander, Charles Cox, Zamin Iqbal, Matthew R. Nelson, and Gil McVean. 2015. “Improved
506 Genome Inference in the MHC Using a Population Reference Graph.” *Nature Genetics* 47 (6): 682–
507 88.
- 508 Dolle, Dirk D., Zhicheng Liu, Matthew Cotten, Jared T. Simpson, Zamin Iqbal, Richard Durbin, Shane A.
509 McCarthy, and Thomas M. Keane. 2017. “Using Reference-Free Compressed Data Structures to
510 Analyze Sequencing Reads from Thousands of Human Genomes.” *Genome Research* 27 (2): 300–
511 309.
- 512 Feuk, Lars, Andrew R. Carson, and Stephen W. Scherer. 2006. “Structural Variation in the Human
513 Genome.” *Nature Reviews Genetics*. <https://doi.org/10.1038/nrg1767>.
- 514 Gardy, Jennifer, Nicholas J. Loman, and Andrew Rambaut. 2015. “Real-Time Digital Pathogen
515 Surveillance — the Time Is Now.” *Genome Biology*. <https://doi.org/10.1186/s13059-015-0726-x>.
- 516 Ghaffaari, Ali, and Tobias Marschall. 2019. “Fully-Sensitive Seed Finding in Sequence Graphs Using a
517 Hybrid Index.” *Bioinformatics* 35 (14): i81–89.
- 518 Hung, Che-Lun, Tzu-Hung Hsu, Hsiao-Hsi Wang, and Chun-Yuan Lin. 2018. “A GPU-Based Bit-Parallel
519 Multiple Pattern Matching Algorithm.” *2018 IEEE 20th International Conference on High
520 Performance Computing and Communications; IEEE 16th International Conference on Smart City;
521 IEEE 4th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*.
522 <https://doi.org/10.1109/hpcc/smartycity/dss.2018.00205>.
- 523 Kaye, Alice. 2016. Methods for the graphical representation of genomic sequence data. USPTO
524 20160342737:A1. *US Patent*, filed May 20, 2016, and issued November 24, 2016.
525 <https://patentimages.storage.googleapis.com/0e/37/3d/b2218ac644a833/US20160342737A1.pdf>.
- 526 Kehrer, Birte, Kathrin Trappe, Manuel Holtgrewe, and Knut Reinert. 2014. “Genome Alignment with
527 Graph Data Structures: A Comparison.” *BMC Bioinformatics* 15 (1): 99.
- 528 Kim, Daehwan, Li Song, Florian P. Breitwieser, and Steven L. Salzberg. 2016. “Centrifuge: Rapid and
529 Sensitive Classification of Metagenomic Sequences.” *Genome Research* 26 (12): 1721–29.
- 530 Klijn, Christiaan, Steffen Durinck, Eric W. Stawiski, Peter M. Haverty, Zhaoshi Jiang, Hanbin Liu,
531 Jeremiah Degenhardt, et al. 2015. “A Comprehensive Transcriptional Portrait of Human Cancer Cell
532 Lines.” *Nature Biotechnology* 33 (3): 306–12.
- 533 Labeit, Julian, Julian Shun, and Guy E. Blelloch. 2017. “Parallel Lightweight Wavelet Tree, Suffix Array
534 and FM-Index Construction.” *Journal of Discrete Algorithms* 43 (March): 2–17.
- 535 Levy-Sakin, Michal, Steven Pastor, Yulia Mostovoy, Le Li, Alden K. Y. Leung, Jennifer McCaffrey,
536 Eleanor Young, et al. 2019. “Genome Maps across 26 Human Populations Reveal Population-
537 Specific Patterns of Structural Variation.” *Nature Communications* 10 (1): 1025.
- 538 Li, Yilong, Nicola D. Roberts, Jeremiah A. Wala, Ofer Shapira, Steven E. Schumacher, Kiran Kumar,
539 Ekta Khurana, et al. 2020. “Patterns of Somatic Structural Variation in Human Cancer Genomes.”
540 *Nature* 578 (7793): 112–21.
- 541 MacDonald, Jeffrey R., Robert Ziman, Ryan K. C. Yuen, Lars Feuk, and Stephen W. Scherer. 2014. “The
542 Database of Genomic Variants: A Curated Collection of Structural Variation in the Human
543 Genome.” *Nucleic Acids Research* 42 (Database issue): D986–92.
- 544 Nielsen, Rasmus, Joshua S. Paul, Anders Albrechtsen, and Yun S. Song. 2011. “Genotype and SNP
545 Calling from next-Generation Sequencing Data.” *Nature Reviews Genetics*.
546 <https://doi.org/10.1038/nrg2986>.
- 547 Olney, Kimberly C., Sarah M. Brotman, Valeria Valverde-Vesling, Jocelyn Andrews, and Melissa A.

- 548 Wilson. n.d. “Aligning RNA-Seq Reads to a Sex Chromosome Complement Informed Reference
549 Genome Increases Ability to Detect Sex Differences in Gene Expression.”
550 <https://doi.org/10.1101/668376>.
- 551 Paten, Benedict, Adam Novak, and David Haussler. 2014. “Mapping to a Reference Genome Structure.”
552 *ArXiv E-Prints*, 1–26.
- 553 Paten, Benedict, Adam M. Novak, Jordan M. Eizenga, and Erik Garrison. 2017. “Genome Graphs and the
554 Evolution of Genome Inference.” *Genome Research* 27 (5): 665–76.
- 555 Patro, Rob, Stephen M. Mount, and Carl Kingsford. 2014. “Sailfish Enables Alignment-Free Isoform
556 Quantification from RNA-Seq Reads Using Lightweight Algorithms.” *Nature Biotechnology* 32 (5):
557 462–64.
- 558 Pedersen, Brent S., and Aaron R. Quinlan. 2017. “Who’s Who? Detecting and Resolving Sample
559 Anomalies in Human DNA Sequencing Studies with Peddy.” *American Journal of Human Genetics*
560 100 (3): 406–13.
- 561 Poplin, Ryan, Pi-Chuan Chang, David Alexander, Scott Schwartz, Thomas Colthurst, Alexander Ku, Dan
562 Newburger, et al. 2018. “A Universal SNP and Small-Indel Variant Caller Using Deep Neural
563 Networks.” *Nature Biotechnology* 36 (10): 983–87.
- 564 Shajii, Ariya, Deniz Yorukoglu, Yun William Yu, and Bonnie Berger. 2016. “Fast Genotyping of Known
565 SNPs through Approximate K-Mer Matching.” *Bioinformatics* 32 (17): i538–44.
- 566 Shen, Shen, and Kidd. 2020. “Rapid, Paralog-Sensitive CNV Analysis of 2457 Human Genomes Using
567 QuicK-mer2.” *Genes*. <https://doi.org/10.3390/genes11020141>.
- 568 Sherman, Rachel M., Juliet Forman, Valentin Antonescu, Daniela Puiu, Michelle Daya, Nicholas Rafaels,
569 Meher Preethi Boorgula, et al. 2019. “Assembly of a Pan-Genome from Deep Sequencing of 910
570 Humans of African Descent.” *Nature Genetics* 51 (1): 30–35.
- 571 Sudmant, Peter H., Tobias Rausch, Eugene J. Gardner, Robert E. Handsaker, Alexej Abyzov, John
572 Huddleston, Yan Zhang, et al. 2015. “An Integrated Map of Structural Variation in 2,504 Human
573 Genomes.” *Nature* 526 (7571): 75–81.
- 574 Sun, Chen, and Paul Medvedev. 2019. “Toward Fast and Accurate SNP Genotyping from Whole Genome
575 Sequencing Data for Bedside Diagnostics.” *Bioinformatics* 35 (3): 415–20.
- 576 Trost, Brett, Susan Walker, Zhuozhi Wang, Bhooma Thiruvahindrapuram, Jeffrey R. MacDonald, Wilson
577 W. L. Sung, Sergio L. Pereira, et al. 2018. “A Comprehensive Workflow for Read Depth-Based
578 Identification of Copy-Number Variation from Whole-Genome Sequence Data.” *American Journal*
579 *of Human Genetics* 102 (1): 142–55.
- 580 Webster, Timothy H., Madeline Couse, Bruno M. Grande, Eric Karlins, Tanya N. Phung, Phillip A.
581 Richmond, Whitney Whitford, and Melissa A. Wilson. 2019. “Identifying, Understanding, and
582 Correcting Technical Artifacts on the Sex Chromosomes in next-Generation Sequencing Data.”
583 *GigaScience* 8 (7). <https://doi.org/10.1093/gigascience/giz074>.
- 584 Wood, Derrick E., Jennifer Lu, and Ben Langmead. 2019. “Improved Metagenomic Analysis with Kraken
585 2.” *Genome Biology* 20 (1): 257.
- 586 Xia, Yuchao, Yun Liu, Minghua Deng, and Ruibin Xi. 2019. “Detecting Virus Integration Sites Based on
587 Multiple Related Sequencing Data by VirTect.” *BMC Medical Genomics* 12 (Suppl 1): 19.
- 588 Yang, Xiaofei, Wan-Ping Lee, Kai Ye, and Charles Lee. 2019. “One Reference Genome Is Not Enough.”
589 *Genome Biology* 20 (1): 104.

590

591

592