

VariantStore: A Large-Scale Genomic Variant Search Index

Prashant Pandey¹, Yinjie Gao¹, and Carl Kingsford^{*1}

¹Computational Biology Department, School of Computer Science, Carnegie Mellon University,
5000 Forbes Ave., Pittsburgh, PA

December 24, 2019

Abstract

The ability to efficiently query genomic variation data from thousands of samples is critical to achieve the full potential of many medical and scientific applications such as personalized medicine. We present VariantStore, a system for efficiently indexing and searching millions of genomic variants across thousands of samples. We show the scalability of VariantStore by indexing genomic variants from the TCGA-BRCA project containing 8640 samples and 5M variants in ≈ 4 Hrs and the 1000 genomes project containing 2500 samples and 924M variants in ≈ 3 Hrs. Querying for variants in a gene takes between 2 milliseconds to 3 seconds using memory only $\approx 10\%$ of the size of the full representation. As a baseline, VariantStore outperformed VG toolkit by $3\times$ in terms of memory-usage and construction time and uses 25% less disk space although VG toolkit does not support variant queries.

Advanced sequencing technology and computing resources have led to large-scale genomic sequencing efforts producing genomic variation data from thousands of samples, such as the 1000 Genomes project [1–3], GTEx [4], and The Cancer Genome Atlas (TCGA) [5]. Analysis of genomic variants combined with phenotypic information of samples promises to improve applications such as personalized medicine, population-level disease analysis, and cancer remission rate prediction. Although numerous studies [6–12] have been performed over the past decade involving genomic variation, the ability to scale these studies to large-scale data available today and in the near future is still limited.

Contrary to the scale of the size of genomic variation data available there is no large-scale indexing and

*To whom correspondence should be addressed: carlk@cs.cmu.edu

24 querying system that can support efficient variant queries. Traditional database solutions, such as SQL and
25 NoSQL, have proven prohibitively slow to store and query collections of variants [13–15].

26 On an individual sample, the typical result of sequencing, alignment, and variant calling is a collection of
27 millions of sample-specific variants. A variant is identified by the position in the chromosome where it
28 occurs, an alternative sequence, a list of samples that contain the variant, and phasing information. The
29 standard file format to report these variants is the variant call file (VCF) [16].

30 A common task is to identify all the samples with a given pattern of variants or identify all samples that
31 have variants in a given gene. These tasks are translated into variant queries which require finding variants
32 or samples with variants between two positions in a chromosome. Positions can be specified in either the
33 reference coordinate or a sample coordinate depending on the query because variants can appear at different
34 coordinates in a sample sequence compared to the reference sequence.

35 To effectively use variant information from many samples, medical and scientific applications must often
36 answer many instances of one or more of the following types of queries:

- 37 1. Find the closest variant to position X for all samples in the reference coordinates.
- 38 2. Find the sequence between positions X and Y for sample S in the reference coordinates.
- 39 3. Find the sequence between positions X and Y for sample S in the sample coordinates.
- 40 4. Find all variants between positions X and Y for sample S in the reference coordinates.
- 41 5. Find all variants between positions X and Y for sample S in the sample coordinates.
- 42 6. Find all variants between positions X and Y for all samples in the reference coordinates.

43 These queries can be further used as building blocks for more complicated queries, such as finding samples
44 with more than a given number of variants between two positions or count the number of variants for each
45 sample between two positions.

46 Genomic variation data are also used for read alignment in order to avoid mapping biases that arise when
47 mapping reads to a single high-quality reference sequence [17–21]. VG toolkit [19] is the most widely used
48 tool to represent genomic variation. It encodes genomic variants from multiple samples in a graph, called
49 a variation graph. A variation graph is a sequence graph where each node represents a sequence and a set
50 of nodes through the graph, known as a path, embeds the complete sequence corresponding to the reference

51 or a sample. Each node on a path is assigned a position corresponding to the location of the node in the
52 coordinate of the reference or a sample sequence the path embeds. Reads are aligned to sequences in the
53 variation graph by following paths in the graph.

54 The variation graph representation in VG toolkit is designed to optimize read alignment and uses sequence-
55 based indexes for alignment. It can not be directly used for variant queries that require an index based on
56 the position of variants in multiple sequence coordinates. VG toolkit stores each sample path as a list of
57 nodes in the graph and maintains a separate index corresponding to the coordinates of the reference and
58 samples. Storing a separate list of nodes for each sequence impedes the scalability of the representation
59 for storing variation from thousands of samples. Moreover, variants are often shared among samples, so
60 storing a list of nodes for each sample path introduces redundancy in the representation. Finally, the VG
61 toolkit representation does not store phasing information contained in VCF files, which is required in many
62 analyses.

63 We present VariantStore, a system for efficiently indexing and querying genomic information (genomic
64 variants and phasing information) from thousands of samples containing millions of variants. VariantStore
65 builds a variant-oriented index by mapping each variant to the list of samples that contain the variant. Vari-
66 ants are indexed based on the positions where they occur. It supports querying variants occurring between
67 two positions across a chromosome based on the reference or a sample coordinate. VariantStore can scale
68 to tens of millions of variants and thousands of samples and can efficiently scale out-of-RAM to storage
69 devices to enable memory-efficient construction and query.

70 We encode genomic variation in a directed, acyclic variation graph and build a position index (a mapping
71 of node positions to node identifiers) on the graph to quickly access a node in the graph corresponding to a
72 queried position. Each node in the variation graph corresponds to a variant and stores a list of samples that
73 contain the variant along with the position of the variant in the coordinate of those samples. The inverted
74 index design allows one to quickly find all the samples and positions in sample coordinates corresponding to
75 a variant. It also avoids redundancy that otherwise arises in maintaining individual variant indexes for each
76 sample coordinate and scales well in practice when the number of samples grows beyond a few thousand.

77 To perform index construction and query efficiently in terms of memory, we partition the variation graph
78 into small chunks (usually a few MBs in size) based on the reference coordinates and store variation graph

79 nodes in these chunks. During construction, an active chunk is always maintained in memory in which new
80 nodes are added, and once it reaches its capacity we compress and serialize it to disk and create a new active
81 chunk. The nodes in and across chunks are ordered based on the reference coordinate since they are created
82 based on variants in the VCF file which are themselves ordered by the reference coordinate. During a query,
83 we only load the chunks in memory that contain the nodes corresponding to the query range.

84 To efficiently scale to thousands of coordinate systems (or samples), we maintain the position index only
85 on the reference coordinate. The position index maps positions in the reference sequence where there is a
86 variant to nodes corresponding to those variants in the variation graph. To lookup a position using a sample's
87 coordinate system, we first lookup the node corresponding to the position in the reference coordinate. We
88 then traverse the sample path from the node in the graph to determine the appropriate node in the sample
89 coordinate. A node with a given position in a sample coordinate is often close to the node in the reference
90 coordinate with the same position.

91 **Results**

92 Our evaluation of VariantStore is based on four parameters: construction time, query throughput, disk space,
93 and peak memory usage.

94 To calibrate our performance numbers, we compare the construction performance and disk space of Vari-
95 antStore against the variation graph representation in VG toolkit [22]. VG toolkit provides a succinct rep-
96 resentation of the variation graph and enables mapping reads to variation graph at the scale of the human
97 genome. VG toolkit is not designed for answering variant queries and therefore is not optimized for it.

98 **Data.** We use 1000 Genomes Phase 3 data [23] and three of the biggest projects from TCGA in terms of the
99 number of samples, Ovarian Cancer (OV), Lung Adenocarcinoma (LUAD), and Breast Invasive Carcinoma
100 (BRCA), for our evaluation. 1000 Genomes data contains more variants compared to the TCGA data but
101 TCGA data contains more samples. 1000 Genomes data contains a separate VCF file for each chromosome
102 containing variants from thousands of samples. The number of samples in each file is $\approx 2.5K$. The variants
103 in 1000 Genomes project are based on GRCh37 reference genome. The TCGA data contains a separate
104 VCF file for each sample. The OV, LUAD, and BRCA projects contain 2436, 2680, and 4319 VCF files
105 containing both normal and tumor variants respectively. For each project in TCGA, we first merged VCF

| System | Time | Disk space | Peak RAM | Peak RAM Agg. |
|--------------|----------------|--------------|----------|---------------|
| Dataset | | 1000 Genomes | | |
| VariantStore | 3 Hrs 25 mins | 41 GB | 8.8 GB | 153 GB |
| VG-toolkit | 11 Hrs 10 mins | 50 GB | 37 GB | 450 GB |
| Dataset | | TCGA (OV) | | |
| VariantStore | 1 Hr 5 mins | 3.4 GB | 1.1 GB | 17.45 GB |
| VG-toolkit | | 11 GB* | | |
| Dataset | | TCGA (LUAD) | | |
| VariantStore | 1 Hr 20 mins | 3.5 GB | 2.3 GB | 36.05 GB |
| VG-toolkit | | 12 GB* | | |
| Dataset | | TCGA (BRCA) | | |
| VariantStore | 4 Hrs 36 mins | 4.2 GB | 3.2 GB | 53.21 GB |
| VG-toolkit | | 14 GB* | | |

Table 1: Time, space, peak RAM, and peak RAM (aggregate) to construct a variation graph and position index on the 1000 Genomes and TCGA (OV, LUAD, and BRCA) data using VariantStore and VG toolkit. We constructed all 24 chromosomes (1 – 22 and X and Y) in parallel. The time and peak RAM reported is for the biggest chromosome (usually chromosome 1 or 2). The space reported is the total space on disk for all 24 chromosomes. The peak RAM (aggregate) is the aggregate peak RAM for all 24 processes. *Space for the XG index that does not contain any path information.

106 files using the BCF tool merge command [24] and created a separate VCF file for each chromosome. The
 107 variants in the TCGA project are based on GRCh38 reference genome.

108 **Index construction.** The total time taken to construct the variation graph representation and index includes
 109 the time taken to read and parse variants from the VCF file, construct the variation graph representation and
 110 indexes, and serialize the final representation to disk. For VariantStore, the reported time includes the time
 111 to create and serialize the position index. The space reported for VariantStore is the sum of the space of the
 112 variation graph representation and position index.

113 For VG toolkit, creating a variation graph representation with multiple coordinate systems (or sample path
 114 annotations) requires creating two indexes, XG and GBWT index. The XG index is a succinct representa-
 115 tion of the variation graph without path annotations that allows memory- and time-efficient random access
 116 operations on large graphs. The GBWT (graph BWT) is a substring index for storing sample paths in the
 117 variation graph. We first create the variation graph representation using the “construct” command including
 118 all sample path annotations. We then create the XG index and GBWT index from the variation graph rep-
 119 resentation to create an index with all sample path annotations in the variation graph. For VG toolkit, the
 120 reported time includes the time to create and serialize the XG and GBWT indexes. The space reported for

121 VG toolkit is the sum of the space of the XG and GBWT indexes. VG toolkit could not build GBWT index
122 on TCGA data even after running for more than a day. We only report space for the XG index (which does
123 not contain any sample path annotations) for TCGA data.

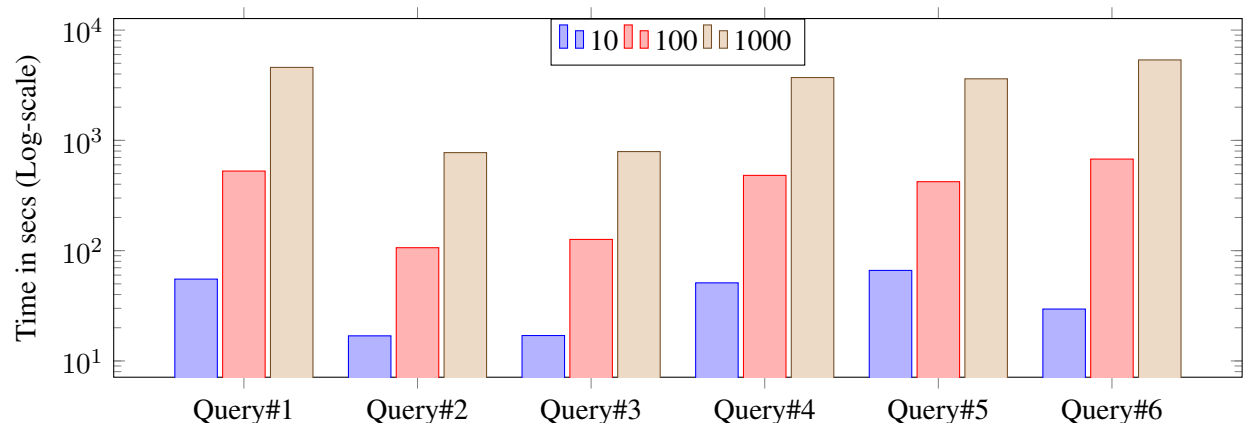
124 For both VariantStore and VG toolkit, we created 24 separate indexes, one each for chromosomes 1 – 22,
125 and X and Y. Each of these indexes were constructed in parallel as a separate process. We report the time
126 taken for construction as the time taken by the process that finishes last. For disk space, we report the total
127 space taken by all 24 indexes. For peak memory usage, we report the highest individual and aggregate peak
128 RAM usage for all processes.

129 The performance of VariantStore and VG toolkit for constructing the index on the 1000 Genomes and TCGA
130 data is shown in Table 1. VariantStore is $3\times$ faster, takes 25% less disk space, and $3\times$ less peak RAM than
131 VG toolkit. For the TCGA data, VG toolkit could not build GBWT index embedding all sample paths.
132 However, even the space needed for the XG index (not embedding sample paths) is $\approx 3.3\times$ larger than the
133 VariantStore representation containing all sample paths.

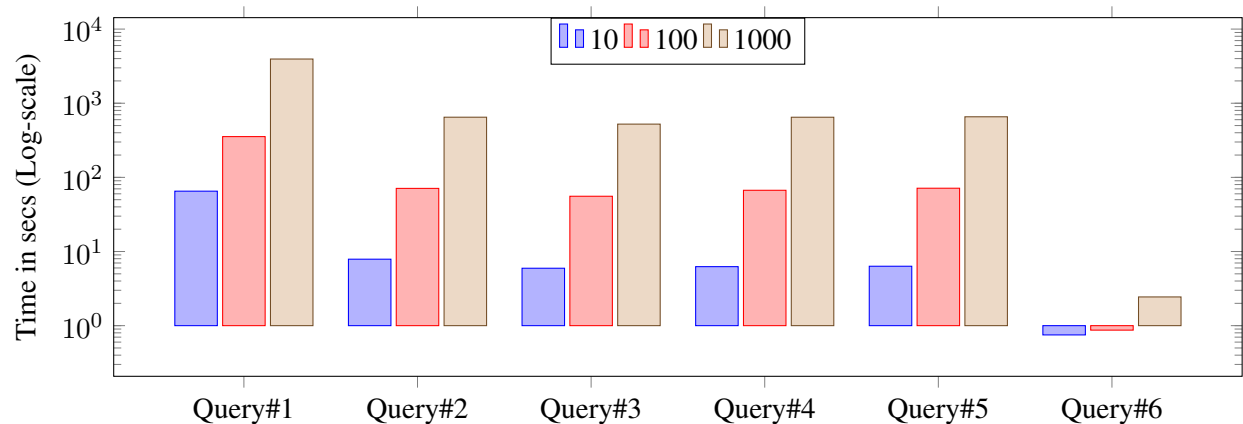
134 **Query throughput.** We measured the query throughput for all six queries mentioned above. To show the
135 robustness of query efficiency on different data we evaluate on two different chromosome indexes from two
136 different projects. We chose Chromosome 2 which is one of the bigger chromosomes and Chromosome 22
137 which is one of the smaller ones. We evaluate query time on chromosome indexes from 1000 Genomes and
138 TCGA LUAD data.

139 To perform queries, we specify a pair of positions in the reference or a sample coordinate system depending
140 on the query type and optionally a sample name. Query parameters, such as the length of the queried region
141 and density of variants in that region, affect the query timing. Therefore, we performed three sets of query
142 benchmarks containing 10, 100, and 1000 queries and report the aggregate time. For each query in the set,
143 we uniformly randomly pick the start position across the full chromosome length. The size of the query
144 range is set to $\approx 42\text{K}$ bases which is approximately the length of a typical gene. Picking multiple query
145 regions uniformly randomly across the chromosome provides a good coverage of different regions across
146 the chromosome.

147 To measure the query throughput, we only load the position index and graph topology in memory and the
148 variation graph representation is kept on disk before performing the query. Keeping the variation graph



(a) Time for 10, 100, and 1000 queries on Chromosome 22 index in VariantStore for 1000 Genomes data.

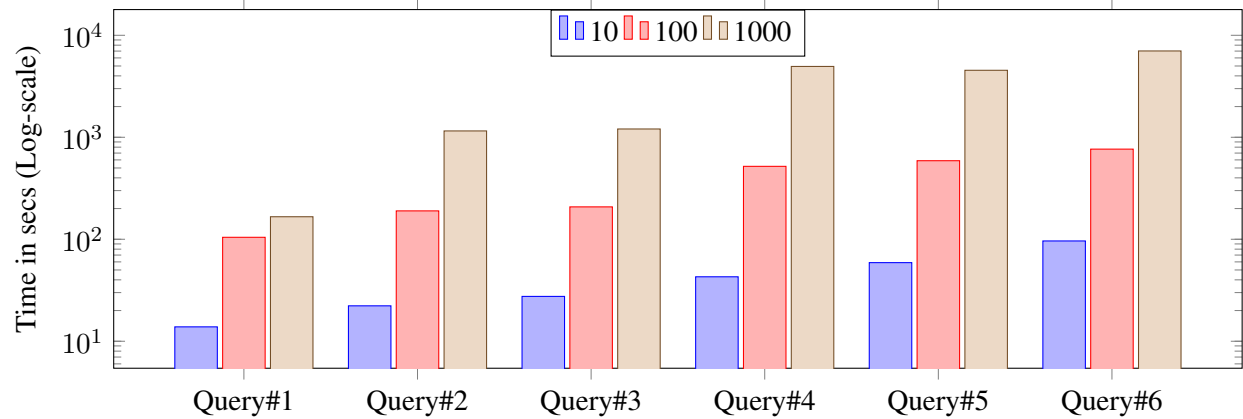


(b) Time for 10, 100, and 1000 queries on Chromosome 22 index in VariantStore for TCGA LUAD data.

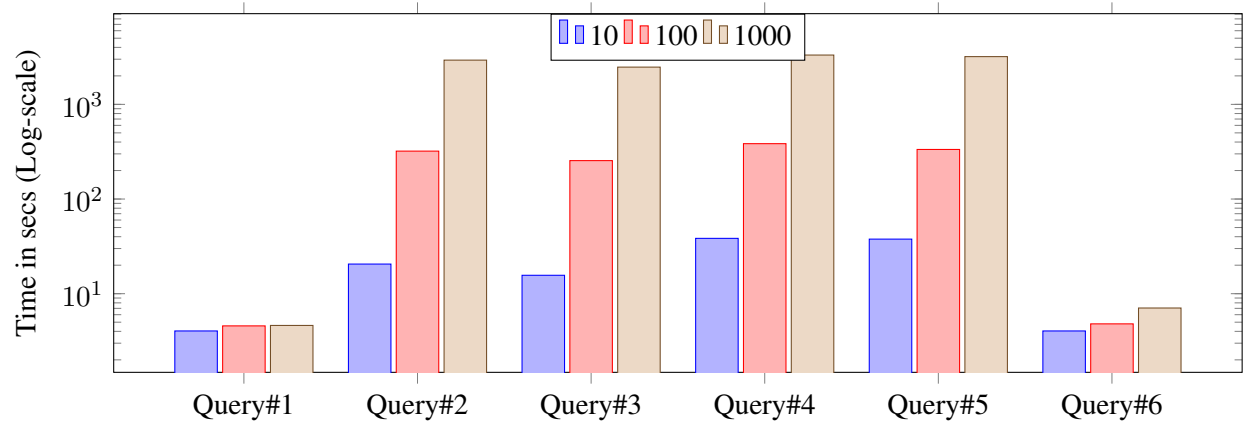
Figure 1: Time reported is the total time taken to execute 10, 100, and 1000 queries. For all queries the query length is fixed to $\approx 42K$. Query#1: Closest variant, Query#2: Seq in Ref coordinate, Query#3: Seq in Sample coordinate, Query#4: Sample vars in Ref coordinate, Query#5: Sample vars in Sample coordinate, Query#6: All vars in Ref coordinate

149 representation on disk keeps the peak memory usage low and all disk accesses are performed during the
 150 query to load appropriate node chunks.

151 The query throughput on 1000 Genomes data is shown in Figures 1a and 2a. For all query types, the
 152 aggregate time taken to execute queries increases linearly with the number of queries. Finding the sequence
 153 corresponding to a sample in a region takes less time compared to finding variants in a region. Finding
 154 the sequence takes less time because it involves traversing the sequence specific path in the region and
 155 reconstructing the sequence. However, finding variants in a given region takes more time because it involves
 156 an exhaustive search of neighbors at each node in the region to determine all the variants that are contained
 157 by a given sample or all samples.



(a) Time for 10, 100, and 1000 queries on Chromosome 2 index in VariantStore for 1000 Genomes data.



(b) Time for 10, 100, and 1000 queries on Chromosome 2 index in VariantStore for TCGA LUAD data.

Figure 2: Time reported is the total time taken to execute 10, 100, and 1000 queries. For all queries the query length is fixed to $\approx 42K$. Query#1: Closest variant, Query#2: Seq in Ref coordinate, Query#3: Seq in Sample coordinate, Query#4: Sample vars in Ref coordinate, Query#5: Sample vars in Sample coordinate, Query#6: All vars in Ref coordinate

158 The query throughput on TCGA data is shown in Figures 1b and 2b. The TCGA data has twice as many
159 samples compared to 1000 Genomes data but there are fewer variants. This makes the variation graph much
160 sparser and queries in the position index become more expensive compared to traversing the graph between
161 two positions. Finding all variants is the fastest query because the variation graph is very sparse and most
162 position ranges are empty. Traversing the graph to find the sequence or variants for a sample takes similar
163 amount of time. Finding the closest variant from a position takes the most amount of time because it involves
164 performing multiple position-index queries to determine the closest variant.

165 For both 1000 Genomes and TCGA LUAD data, finding the closest variant query is faster for Chromosome
166 2 because variants in Chromosome 22 are more dense compared to Chromosome 2 which makes it faster to

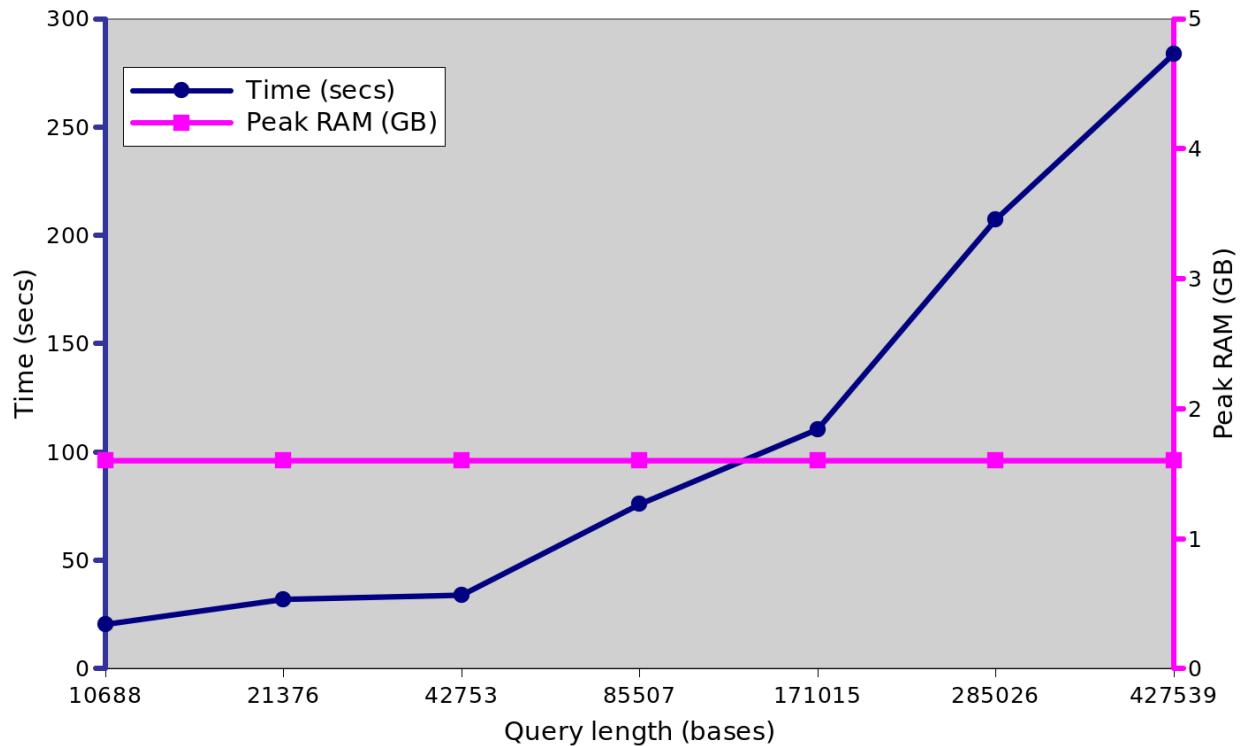


Figure 3: Time (seconds) and peak memory usage (GB) with increasing query length for “Sample sequence in the reference coordinate” (query #2) for 1000 Genomes chromosome 22 index. The time taken increases as the query length increases. But the memory usage remains constant regardless of the query length.

167 locate the closest variant.

168 **The effect of query range on peak memory usage.** We performed another query benchmark to evaluate
169 the effect of size of the position range on the peak memory usage and time. For this benchmark, we chose the
170 “Sample sequence in reference coordinate” query (query #2) because this query involves traversing the full
171 sample path between two positions. We performed sets of 100 queries with increasing size of the position
172 range and record the total time and peak RAM usage. For each query in the set, we uniformly randomly
173 pick the start position across the full chromosome length.

174 Effect of the query range size on peak memory usage and time is shown in Figure 3. The memory usage
175 remains constant regardless of the query length. This is because during a query we access node chunks in
176 sequential order and regardless of the query length only load at most two node chunks in RAM at a time.
177 This keeps the memory usage essentially constant.

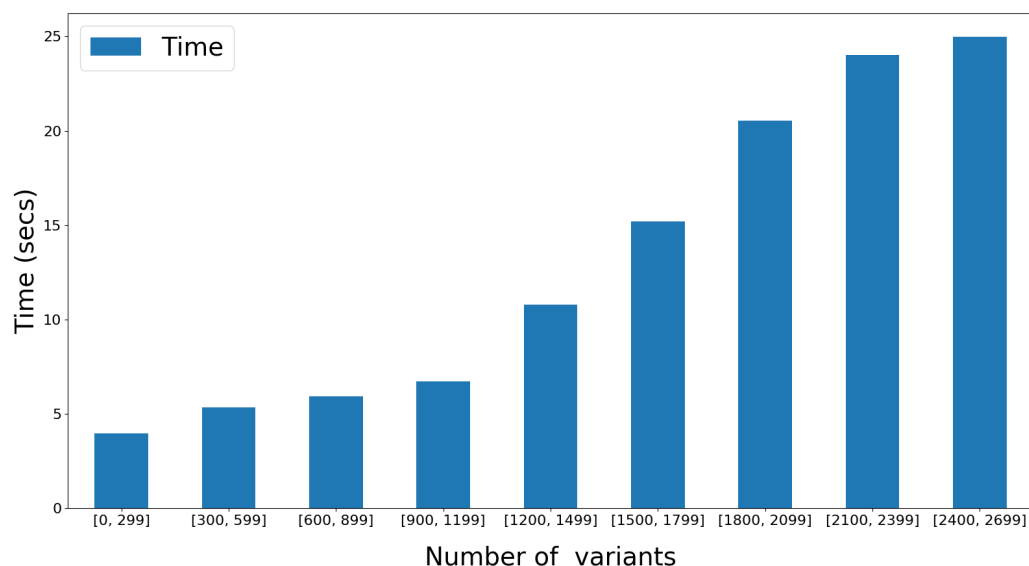


Figure 4: Time (seconds) and the number of variants in the query region for “All variants in the reference coordinate” (query #6) for 1000 Genomes chromosome 22 index. Query times are binned based on the number of variants in the query region and mean time is reported in each bin. The mean time increases as the number of variants increases in the region.

178 **The effect of number of variants on query time.** We also evaluate how the number of variants in the
179 position range affects the query time. For this benchmark, we performed the “All variants in the reference
180 coordinate” query (query #6) because this query involves performing a breadth-first search in the graph to
181 determine all variants in a region and the query time depends on the number of variants in the region. To
182 perform queries on regions with different number of variants, we chose 1000 regions with a fixed size of the
183 position range ($\approx 42\text{K}$ bases) and start position chosen uniformly randomly across the chromosome.

184 Effect of the number of variants in the query region on query time is shown in Figure 4. The query time
185 increases as the number of variants in the queried region increases. This is because when the number of
186 variants in a region is small the graph is sparser and faster to traverse and report all variants.

187 **Experimental hardware.** All experiments on 1000 Genomes data were performed on an Intel Xeon CPU
188 E5-2699A v4 @ 2.40GHz (44 cores and 56MB L3 cache) with 1TB RAM and a 7.3TB HGST HDN728080AL
189 HDD running Ubuntu 18.04.2 LTS (Linux kernel 4.15.0-46-generic) and were run using a single thread.
190 Benchmarks for TCGA data were performed on a cluster machine running AMD Opteron Processor 6220
191 @ 3GHz with 6MB L3 cache. TCGA data was stored on a remote disk and accessed via NFS.

192 **Discussion**

193 We attribute the scalability and efficient index construction and query performance of VariantStore to the
194 variant-oriented index design. VariantStore uses an inverted index from variants to the samples which scales
195 efficiently when multiple samples share a variant which is often seen in genomic variation data. The inverted
196 index design further allows us to build the position index only on the reference sequence and use graph
197 traversal to transform the position in reference coordinates to sample coordinates.

198 All the supported variant queries look for variants in a contiguous region of the chromosome which allows
199 VariantStore to partition the variation graph representation into small chunks based on the position of nodes
200 in the chromosome and sequentially load only the relevant chunks into memory during a query. This makes
201 VariantStore memory-efficient and scale to genomic variation data from hundreds of thousands of samples
202 in future.

203 The variation graph representation in VariantStore is smaller and more efficient to construct than the rep-
204 resentation in state-of-the-art VG toolkit. It can be further used in read alignment as a replacement to the
205 variation graph representation in the VG toolkit.

206 While the current implementation does not support adding new variants or update the reference sequence in
207 an existing VariantStore, this is not a fundamental limitation of the design. In future, we plan to extend the
208 immutable version of VariantStore to support dynamic updates following the LSM-tree design [25].

209 **Methods**

210 **Variation graph**

211 A variation graph (VG) [19] (also defined as a genome graph in Kim et al. [20] and Rakocevic et al. [21]) is
212 a directed, acyclic graph (DAG) $G = (N, E, P)$ that embeds a set of DNA sequences. It comprises of a set
213 of nodes N , a set of direct edges E , and a set of paths P . For DNA sequences, we use the alphabet $\{A, C,$
214 $G, T, N\}$. Each $n_i \in N$ represents a sequence $seq(n_i)$. Edges in the graph connect nodes that are followed
215 on a path. Each sample in the VCF file follows a path through the variation graph. The embedded sequence
216 given by the path is the sample sequence. Given that a variation graph is a directed graph edges can be

217 traversed in only one direction. Although, not applicable to VariantStore, an edge can also be traversed in
218 the reverse direction when the variation graph is used for read alignment [19–21].

219 Each path $p = n_s n_1 \dots n_p n_d$ in the graph is an embedded sequence defined as a sequence of nodes between
220 a source node n_s and a destination node n_d . Nodes on a path are assigned positions based on the coordinate
221 systems of sequences they represent [26]. The position of a node on a path is the sum of the lengths of the
222 sequences represented by nodes traversed to get to the node on the path. For a path $p = n_1 \dots n_p$, position
223 $P(n_p)$ is $\sum_{i=1}^{p-1} |seq(n_i)|$. A node in the graph can appear in multiple paths and therefore can have multiple
224 positions based on different coordinate systems.

225 An initial variation graph is constructed with a single node and no edges using a linear reference sequence
226 and a reference genome coordinate system. Variants are added to the variation graph from one or more VCF
227 files [16]. A variant is encoded by a node in the variation graph that represents the variant sequence and
228 is connected to nodes representing the reference sequence via directed edges. Each variant in the VCF file
229 splits an existing reference sequence node into two (or three in some cases) and joins them via an alternative
230 path corresponding to the variant. For example, a substitution or deletion can cause an existing reference
231 node to be split into three parts (Figure 5). An insertion can cause the reference node to be split into two
232 parts (Figure 5).

233 **Representing multiple coordinate systems in a variation graph**

234 Representing multiple coordinate systems in a variation graph poses challenges that are not present in linear
235 reference genomes. First, a node can appear on multiple paths at a different position on each path. Second,
236 given that a variation graph can contain thousands of paths and coordinate systems it would be non-trivial
237 to maintain a position index to quickly get to a node corresponding to a path and position.

238 Much work has been done in coming up with efficient approaches to handle multiple coordinate sys-
239 tems. Rand et al. [26] introduced the offset-based coordinate system. VG toolkit [22] implemented the
240 multiple coordinate system by explicitly storing a list of node identifiers and node offsets for each path in
241 the variation graph. They store the list of node identifiers as integer vectors using the succinct data structure
242 library (SDSL [27]). We call this an *explicit-path representation*.

243 However, storing the list of nodes on each path explicitly can become a bottleneck as the number of input

| Reference sequence | | CAATTTGCTGATCT | | | | |
|--------------------|----------------|------------------|---------|---------|---------|--------------|
| Position | Reference seq. | Alternative seq. | HG00096 | HG00101 | HG00103 | Variant type |
| 2 | A | G | 0 | 1 | 1 | SUBSTITUTION |
| 2 | AATT | A | 1 | 0 | 0 | DELETION |
| 6 | T | TACG | 0 | 0 | 1 | INSERTION |

Table 2: Variants ordered by the position in the reference genome for three samples (HG00096, HG00101, HG00103). Each variant has the list of samples that contain the variant.

244 paths increases. Moreover, nodes that appear on multiple paths are stored multiple times causing redundancy
245 in storage. For a set of N variants and S samples the space required to store the explicit-path representation
246 is $O(SN)$ since each variant creates a constant number of new nodes in the variation graph.

247 VariantStore

248 We describe how we represent a variation graph in VariantStore and maintain multiple coordinate systems
249 efficiently. We then describe how we build a position index using succinct data structures.

250 The variation graph representation is divided into three components:

- 251 1. Variation graph topology
- 252 2. Sequence buffer
- 253 3. List of variation graph nodes

254 **Variation graph topology.** A variation graph constructed by inserting variants from VCF files often shows
255 high sparsity (the number of edges is close to the number of nodes). For example, the ratio of the number
256 of edges to nodes in the variation graph on 1000 Genomes data [23] is close to 1. Given the sparsity of the
257 graph, we store the topology of the variation graph in a representation optimized for sparse graphs.

258 Our graph representation uses the counting quotient filter (CQF) [28] as the underlying container for storing
259 nodes and their outgoing neighbors. The CQF is a compact representation of a multiset S . Thus a CQF
260 supports inserts and queries. A query for an item x returns the number of instances of x in S . The CQF
261 uses a variable-length encoding scheme to count the multiplicity of items. In our graph representation, we
262 encode the node as the item and id of the outgoing neighbor as the count of the item in the CQF.

263 For nodes with a single outgoing neighbor, we store the node and its neighbor together (without an indirec-
264 tion) so we can access them quickly. For nodes with more than one outgoing neighbor we store outgoing

265 neighbors in separate lists. In the variation graph, most nodes have a single outgoing neighbor and using
266 this compact and optimized representation we achieve cache efficient and fast traversal of the graph.

267 We use the version of the counting quotient filter with no false-positives to map a node to its outgoing
268 edge(s). We store the node id as the key in the CQF and if there is only one outgoing edge we encode the
269 outgoing neighbor id as the count of the key. If there are more than one outgoing edges we use indirection.
270 We maintain a list of vectors where each vector contains a list of outgoing neighbor identifiers corresponding
271 to a node. We store the node id as the key and the offset in the list of vectors (or index of the vector containing
272 the list of outgoing neighbor ids) as the count of the key.

273 **Sequence buffer.** The sequence buffer contains the reference sequence and all variant sequences corre-
274 sponding to each substitution and insertion variant. All sequences are encoded using 3-bit characters in an
275 integer vector from SDSL library [27, 29]. The integer vector initially only contains the reference sequence.
276 Sequences from incoming variants are appended to the integer vector. Once all variants are inserted the
277 integer vector is bit compressed before being written to disk.

278 **List of variation graph nodes.** Each node in the variation graph contains an offset and length. The offset
279 points to the start of the sequence in the sequence buffer and length is the number of nucleotides in the
280 sequence starting from the offset. This uniquely identifies a node sequence in the sequence buffer.

281 At each node we also store a list of sample identifiers that have the variant, position of the node on all those
282 sample paths, and phasing information from the VCF file corresponding to each sample.

283 Our representation of the list of samples is based on two observations. First, multiple samples share a variant
284 and storing a list of sample identifiers for each variant is space inefficient. Instead, we store a bit vector of
285 length equal to the number of samples and set bits corresponding to the present samples in the bit vector.
286 Second, multiple variants share the same set of samples. We define an equivalence relation \sim over the set
287 of variants. Let $E(v)$ denote the function that maps each variant to the set of samples that have the variant.
288 We say that two variants are equivalent (i.e., $v_1 \sim v_2$) if and only if $E(v_1) = E(v_2)$. We refer to the set of
289 samples shared by variants as *sample class*. A unique id is assigned to each sample class and nodes store
290 the sample class id instead of the whole sample class. This scheme is previously been employed by other
291 colored de Bruijn graph representation tools [30–32] for efficiently maintaining a mapping from k -mers (a
292 k -length substring sequence) to the set of samples where k -mers appear.

293 Phasing information is encoded using 3 bits. Position and phasing information corresponding to each sample
294 in the list of samples is stored as tuples. Tuples are stored in the same order as the samples appear in the
295 sample class bit vector. To retrieve the tuple corresponding to a sample a rank operation is performed on the
296 sample bit vector to determine the rank of the sample. Using the rank output, a select operation is performed
297 on the tuple list to determine the tuple corresponding to a sample.

298 Variation graph nodes are stored as protocol buffer objects using Google's open-source protocol buffers
299 library. Every time a new node is created we instantiate a new protocol buffer object in memory. We
300 compress the protocol buffers before writing them to disk and decompress them while reading them back in
301 memory.

302 For a set of N variants and S samples where each variant is shared by P samples on average, each node
303 contains information about P samples and storing $O(N)$ nodes (a constant number of nodes for each variant)
304 the space required to store the variation graph representation in VariantStore is $O(NP)$. When $P = 1$ (i.e.,
305 no two samples share a variant) the space required by the variation graph representation becomes $O(N)$.

306 **Position index** In order to answer variant queries we need an index to quickly locate nodes in the graph
307 corresponding to input positions. These position can be specified in multiple coordinate systems, i.e., in the
308 coordinate system of the reference or a sample.

309 One way to index the variation graph, is to store an ordered mapping from position \rightarrow node identifier. We can
310 perform a binary search in the map to find the position closest to the queried position and the corresponding
311 node id in the graph. However, given that there are multiple coordinate systems in the variation graph we can
312 not create a single mapping with a global ordering. Keeping a separate position index for each coordinate
313 system will require space equal to the explicit-path representation.

314 In VariantStore, we maintain a mapping of positions to node identifiers only for the reference coordinate
315 system. All nodes on the reference sequence path are present in the mapping. If the queried position is in
316 the reference coordinate system we use the mapping to locate the node in the graph. However, if the queried
317 position is in a sample's coordinate system, we first locate the node in the graph corresponding to the same
318 position in the reference coordinate system. Then we perform a local search by traversing the sample path
319 from that node to determine the node corresponding to the position in sample's coordinate system. The local
320 graph search incurs a small one-time cost because sample nodes are rarely far from a reference node and is

321 amortized against future searches in the sample's coordinate system.

322 We create the position index using a bit vector called the position-bv of length equal to the reference se-
323 quence length and a list of node identifiers on the reference path in the increasing order by their reference
324 sequence positions. For every node in the list we set the bit corresponding to the node's position in the
325 position-bv. There is a one-to-one correspondence between every set bit in the position-bv and node posi-
326 tions in the list. We store the position-bv using a bit vector and node list as an integer vector from the SDSL
327 library [27, 29].

328 **Variation graph construction**

329 We construct the variation graph by inserting variants from a VCF file. Each variant has a position in the
330 reference genome, alternative sequence (except in case of a deletion), and a list of samples with phasing
331 information for each sample.

332 Based on the position of the variant we split an existing reference node that contains the sequence at that
333 position in the graph. We update the split nodes on the reference path with new sequence buffer offsets,
334 lengths, and node positions (based on the reference coordinate system). We then append the alternative
335 sequence to the sequence buffer and create an alternative node with the offset and length of the alternative
336 sequence. We then add the list of tuples (position, phasing info) for each sample.

337 We also need to determine the position of the alternative node on the path of each sample that contains the
338 variant. One way to determine the position of the node for each sample would be to backtrack in the graph
339 to determine a previous node that contains a sample variant and the absolute position of that node in the
340 sample's coordinate system. If no node is found with a sample variant we trace all the way back to the
341 source of the graph. We would then traverse the sample path forward up to the new alternative node and
342 compute the position. This backtracking process would need to be performed once for each sample that
343 contains the variant. This would slow down adding a new variant and cause the construction process to not
344 scale well with increasing number of samples.

345 Instead, we construct the variation graph in two phases to avoid the backtracking process. In the first phase,
346 while adding variants we do not update the position of nodes on sample paths. We only maintain the position
347 of nodes on the reference path because that does not require backtracking. In the second phase, we perform

348 a breadth-first traversal of the variation graph starting from the source node and update the position of nodes
349 on sample paths.

350 During the breadth-first traversal we maintain a delta value for each sample in the VCF file. At any node, the
351 *delta value* is the difference between the position of the node in the reference coordinate and the sample's
352 coordinate. During the traversal, we update sample positions for each node based on the current delta value
353 and reference coordinate value. Algorithm 1 gives the pseudocode of the algorithm.

Algorithm 1 Fix sample positions

```
1: for  $i$  in Samples do
2:    $\delta[i] \leftarrow 0$ 
3: for  $node$  in BFS(variation graph) do
4:   if ISREFERENCE( $node$ ) then
5:     for  $neighbor$  in  $node.neighbors$  do
6:       if  $pos[sample] = 0$  then
7:          $pos[sample] \leftarrow pos[ref] + node.len + \delta[sample]$ 
8:       else
9:          $\delta[sample] \leftarrow pos[sample] - (pos[ref] + node.len)$ 
10:  else
11:     $\delta[sample] \leftarrow pos[sample] + pos[ref]$ 
```

354 Position index construction

355 In the position index, we maintain a mapping from positions of nodes on the reference path to corresponding
356 node identifiers in the graph. Node positions are stored in a “position-bv” bit vector of size equal to the
357 length of the reference sequence and node identifiers are stored in a list. To construct the position index, we
358 follow the reference path starting from the source node in the graph and for every node on the path we set
359 the corresponding position bit in the position-bv and add the node identifier to the list. Node identifiers are
360 stored in the order of their position on the reference path.

361 Variant queries

362 A query is performed in two steps. We first perform a predecessor search (largest item smaller than or equal
363 to the queried item) using the queried position in the position index to locate the node n_p with the highest
364 position smaller than or equal to the queried position pos . The predecessor search is implemented using
365 the rank operation on position-bv. For bit vector $B[0, \dots, n]$, $RANK(j)$ returns the number of 1s in prefix

366 $B[0, \dots, j]$ of B . An RRR compressed position-bv supports rank operation in constant time [27, 33]. The
367 rank of pos in position-bv corresponds to the index of the node id in the node list. Figure 6a shows a sample
368 query in the position index.

369 Based on how reference nodes are split while adding variants the sequence starting at pos will either be
370 contained in the node n_p or the next node after n_p on the reference path. All queries are then answered by
371 traversing the graph either by following a specific path (reference or a sample) or a breadth-first traversal
372 and filtering nodes based on query options.

373 If the queried position is based on the reference coordinate system then we can directly use n_p as the start
374 node for graph traversal. However, if the position is based on a sample coordinate then we perform a local
375 search in the graph starting from n_p to determine the start node based on the sample coordinate.

376 **Memory-efficient construction and query**

377 In the variation graph representation, the biggest component in terms of space is the list of variation graph
378 nodes stored as Google protobuf objects. These node objects contain the sequence information and the list
379 of sample positions and phasing information. For 1000 Genomes data, the space required for variation graph
380 nodes is $\approx 87\%$ to 92% of the total space in VariantStore. However, keeping the full list of node objects in
381 memory during construction or query is not necessary and would make these processes memory inefficient.

382 To perform memory efficient construction and query, we store and serialize these nodes in small chunks
383 usually containing $\approx 200\text{K}$ nodes (the number of nodes in a chunk varies based on the data to keep the
384 size to a few MBs). Nodes in and across these chunks are kept in their creation order (which is roughly the
385 breadth-first traversal order). Therefore, during a breadth-first traversal of the graph we only need to load
386 these chunk in sequential order.

387 During construction, we only keep two chunks in memory, the current active chunk and the previous one.
388 All chunks before the previous chunk are written to disk. In the second phase of the construction when we
389 update sample positions and during the position index creation we perform a breadth-first traversal on the
390 graph and load chunks in sequential order.

391 Variant queries involve traversing a path in the graph between a start and an end position or exploring the

392 graph locally around a start position. All these queries require bounded exploration of the graph for which
393 we only need to look into one or a few chunks.

394 To perform queries with a constant memory we only load the position index and variation graph topology
395 in memory and keep the node chunks on disk. We use the index and the graph topology to determine the set
396 of nodes to look at to answer the query. We then load appropriate chunks from disk which contain the start
397 and end nodes in the query range. For queries involving local exploration of the graph we load the chunk
398 containing the start node. During the exploration, we load new chunks lazily as needed. At any time during
399 the query, we only maintain two contiguous chunks in memory.

400 **References**

- 401 [1] 1000 Genomes Project Consortium. A map of human genome variation from population-scale se-
402 quencing. *Nature*, 467(7319):1061, 2010.
- 403 [2] 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human
404 genomes. *Nature*, 491(7422):56, 2012.
- 405 [3] 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*, 526
406 (7571):68, 2015.
- 407 [4] John Lonsdale, Jeffrey Thomas, Mike Salvatore, Rebecca Phillips, Edmund Lo, Saboor Shad, Richard
408 Hasz, Gary Walters, Fernando Garcia, and Nancy Young. The genotype-tissue expression (GTEx)
409 project. *Nature Genetics*, 45(6):580, 2013.
- 410 [5] TCGA: the cancer genome atlas program, 2019. URL [https://www.cancer.gov/about-nci/organization/
411 ccg/research/structural-genomics/tcga](https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga). [Online accessed August 2019].
- 412 [6] Ananyo Choudhury, Michèle Ramsay, Scott Hazelhurst, Shaun Aron, Soraya Bardien, Gerrit Botha,
413 Emile R Chimusa, Alan Christoffels, Junaid Gamiieldien, and Mahjoubeh Sefid-Dashti. Whole-genome
414 sequencing for an enhanced understanding of genetic variation among South Africans. *Nature Com-
415 munications*, 8(1):2062, 2017.
- 416 [7] Sebastian Roskosch, Hákon Jónsson, Eythór Björnsson, Doruk Beyter, Hannes P Eggertsson, Patrick

- 417 Sulem, Kári Stefánsson, Bjarni V Halldórsson, and Birte Kehr. PopDel identifies medium-size dele-
418 tions jointly in tens of thousands of genomes. *bioRxiv*, page 740225, 2019.
- 419 [8] Cristian Groza, Tony Kwan, Nicole Soranzo, Tomi Pastinen, and Guillaume Bourque. Personalized
420 and graph genomes reveal missing signal in epigenomic data. *bioRxiv*, page 457101, 2019.
- 421 [9] Frank W Albert and Leonid Kruglyak. The role of regulatory variation in complex traits and disease.
422 *Nature Reviews Genetics*, 16(4):197, 2015.
- 423 [10] Karen Eilbeck, Aaron Quinlan, and Mark Yandell. Settling the score: variant prioritization and
424 mendelian disease. *Nature Reviews Genetics*, 18(10):599, 2017.
- 425 [11] Claudia MB Carvalho and James R Lupski. Mechanisms underlying structural variant formation in
426 genomic disorders. *Nature Reviews Genetics*, 17(4):224, 2016.
- 427 [12] Jerome Kelleher, Yan Wong, Anthony W Wohns, Chaimaa Fadil, Patrick K Albers, and Gil McVean.
428 Inferring whole-genome histories in large population datasets. *Nature Genetics*, 51(9):1330–1338,
429 2019.
- 430 [13] Anthony J Brookes and Peter N Robinson. Human genotype–phenotype databases: aims, challenges
431 and opportunities. *Nature Reviews Genetics*, 16(12):702, 2015.
- 432 [14] Joachim Kutzera and Patrick May. Variant-DB: A tool for efficiently exploring millions of human
433 genetic variants and their annotations. In *International Conference on Data Integration in the Life*
434 *Sciences*, pages 22–28. Springer, 2017.
- 435 [15] Geert Vandeweyer, Lut Van Laer, Bart Loeys, Tim Van den Bulcke, and R Frank Kooy. VariantDB: a
436 flexible annotation and filtering portal for next generation sequencing data. *Genome Medicine*, 6(10):
437 74, 2014.
- 438 [16] The variant call format (VCF) version 4.1 specification, 2019. URL [https://samtools.github.io/](https://samtools.github.io/hts-specs/VCFv4.1.pdf)
439 [hts-specs/VCFv4.1.pdf](https://samtools.github.io/hts-specs/VCFv4.1.pdf). [Online accessed March 2019].
- 440 [17] Alexander Dilthey, Charles Cox, Zamin Iqbal, Matthew R Nelson, and Gil McVean. Improved genome
441 inference in the MHC using a population reference graph. *Nature Genetics*, 47(6):682, 2015.
- 442 [18] Hannes P Eggertsson, Hakon Jonsson, Snaedis Kristmundsdottir, Eirikur Hjartarson, Birte Kehr,

- 443 Gisli Masson, Florian Zink, Kristjan E Hjorleifsson, Aslaug Jonasdottir, and Adalbjorg Jonasdottir.
444 GraphTyper enables population-scale genotyping using pangenome graphs. *Nature Genetics*, 49(11):
445 1654, 2017.
- 446 [19] Erik Garrison, Jouni Sirén, Adam M Novak, Glenn Hickey, Jordan M Eizenga, Eric T Dawson, William
447 Jones, Shilpa Garg, Charles Markello, Michael F Lin, Benedict Paten, and Richard Durbin. Varia-
448 tion graph toolkit improves read mapping by representing genetic variation in the reference. *Nature*
449 *Biotechnology*, 36:875–879, 2018.
- 450 [20] Daehwan Kim, Joseph M Paggi, Chanhee Park, Christopher Bennett, and Steven L Salzberg. Graph-
451 based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology*,
452 37(8):907–915, 2019.
- 453 [21] Goran Rakocevic, Vladimir Semenyuk, Wan-Ping Lee, James Spencer, John Browning, Ivan J Johnson,
454 Vladan Arsenijevic, Jelena Nadj, Kaushik Ghose, and Maria C Suciuc. Fast and accurate genomic
455 analyses using genome graphs. *Nature Genetics*, 51:354–362, 2019.
- 456 [22] Variation graph toolkit, 2019. URL <https://github.com/vgteam/vg>. [Online accessed March 2019].
- 457 [23] IGSR: the international genome sample resource, 2019. URL <http://www.internationalgenome.org/>
458 home. [Online accessed March 2019].
- 459 [24] BCF toolkit, 2019. URL <https://samtools.github.io/bcftools/bcftools.html>. [Online accessed March
460 2019].
- 461 [25] Patrick O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth O’Neil. The log-structured merge-tree
462 (LSM-tree). *Acta Informatica*, 33(4):351–385, 1996.
- 463 [26] Knut D Rand, Ivar Grytten, Alexander J Nederbragt, Geir O Storvik, Ingrid K Glad, and Geir K
464 Sandve. Coordinates and intervals in graph-based reference genomes. *BMC Bioinformatics*, 18(1):
465 263, 2017.
- 466 [27] Simon Gog, Timo Beller, Alistair Moffat, and Matthias Petri. From theory to practice: Plug and play
467 with succinct data structures. In *13th International Symposium on Experimental Algorithms, (SEA*
468 *2014)*, pages 326–337, 2014.

- 469 [28] Prashant Pandey, Michael A Bender, Rob Johnson, and Rob Patro. A general-purpose counting filter:
470 Making every bit count. In *Proceedings of the 2017 ACM International Conference on Management*
471 *of Data*, pages 775–787. ACM, 2017.
- 472 [29] SDSL: succinct data structure library, 2019. URL <https://github.com/simongog/sdsl-sdslite>. [Online
473 accessed March 2019].
- 474 [30] Fatemeh Almodaresi, Prashant Pandey, and Rob Patro. Rainbowfish: a succinct colored de Bruijn
475 graph representation. In *17th International Workshop on Algorithms in Bioinformatics (WABI 2017)*.
476 Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- 477 [31] Prashant Pandey, Fatemeh Almodaresi, Michael A Bender, Michael Ferdman, Rob Johnson, and Rob
478 Patro. Mantis: A fast, small, and exact large-scale sequence-search index. *Cell Systems*, 7(2):201–207,
479 2018.
- 480 [32] Fatemeh Almodaresi, Prashant Pandey, Michael Ferdman, Rob Johnson, and Rob Patro. An efficient,
481 scalable and exact representation of high-dimensional color information enabled via de Bruijn graph
482 search. In *International Conference on Research in Computational Molecular Biology*, pages 1–18.
483 Springer, 2019.
- 484 [33] Rajeev Raman, Venkatesh Raman, and Srinivasa Rao Satti. Succinct indexable dictionaries with appli-
485 cations to encoding k-ary trees, prefix sums and multisets. *ACM Transactions on Algorithms (TALG)*,
486 3(4):43, 2007.

487 **Data availability**

488 Code and data are available at <https://github.com/Kingsford-Group/variantstore>

489 **Acknowledgements**

490 The results published here are in whole or part based upon data generated by The Cancer Genome Atlas
491 (dbGaP accession phs000178) managed by the NCI and NHGRI. Information about TCGA can be found
492 at <http://cancergenome.nih.gov> The 1000 Genomes data used for the analyses described in this manuscript

493 were obtained from <https://www.internationalgenome.org/>. This research is funded in part by the Gordon
494 and Betty Moore Foundation's Data-Driven Discovery Initiative through Grant GBMF4554 to C.K. and by
495 the US National Institutes of Health (R01GM122935). We would also like to thank Shawn Baker for many
496 helpful discussions, and Guillaume Marçais and Yutong Qiu for comments on the manuscript.

497 **Competing Interests**

498 C.K. is a co-founder of Ocean Genomics, Inc.

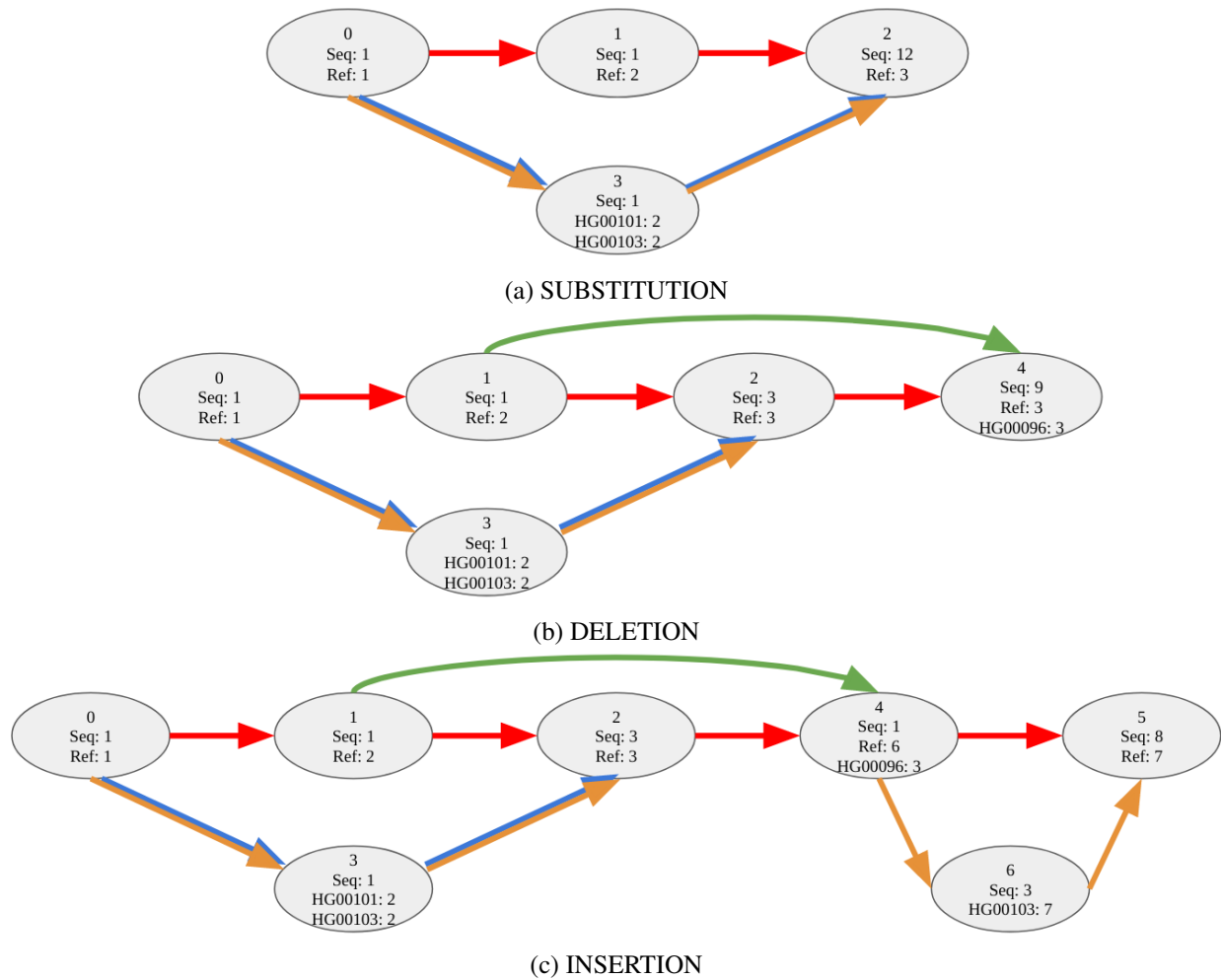


Figure 5: A variation graph with three input samples (HG00096, HG00101, HG00103) showing the encoding of substitutions, insertions, and deletions as stated in Table 2. Fig. (a) shows the substitution, (b) shows the deletion, and (c) shows the insertion. Edges are colored (or multi-colored) to show the path taken by reference and samples through the graph. (Ref: red, HG00096: green, HG00101: blue, HG00103: brown). Samples with no variant at a node follow the reference path, e.g., sample HG00096 will follow the reference path between nodes $0 \rightarrow 1$ and $4 \rightarrow 5$. Each node contains node id, the length of the sequence it represents, and a list of samples and their positions.

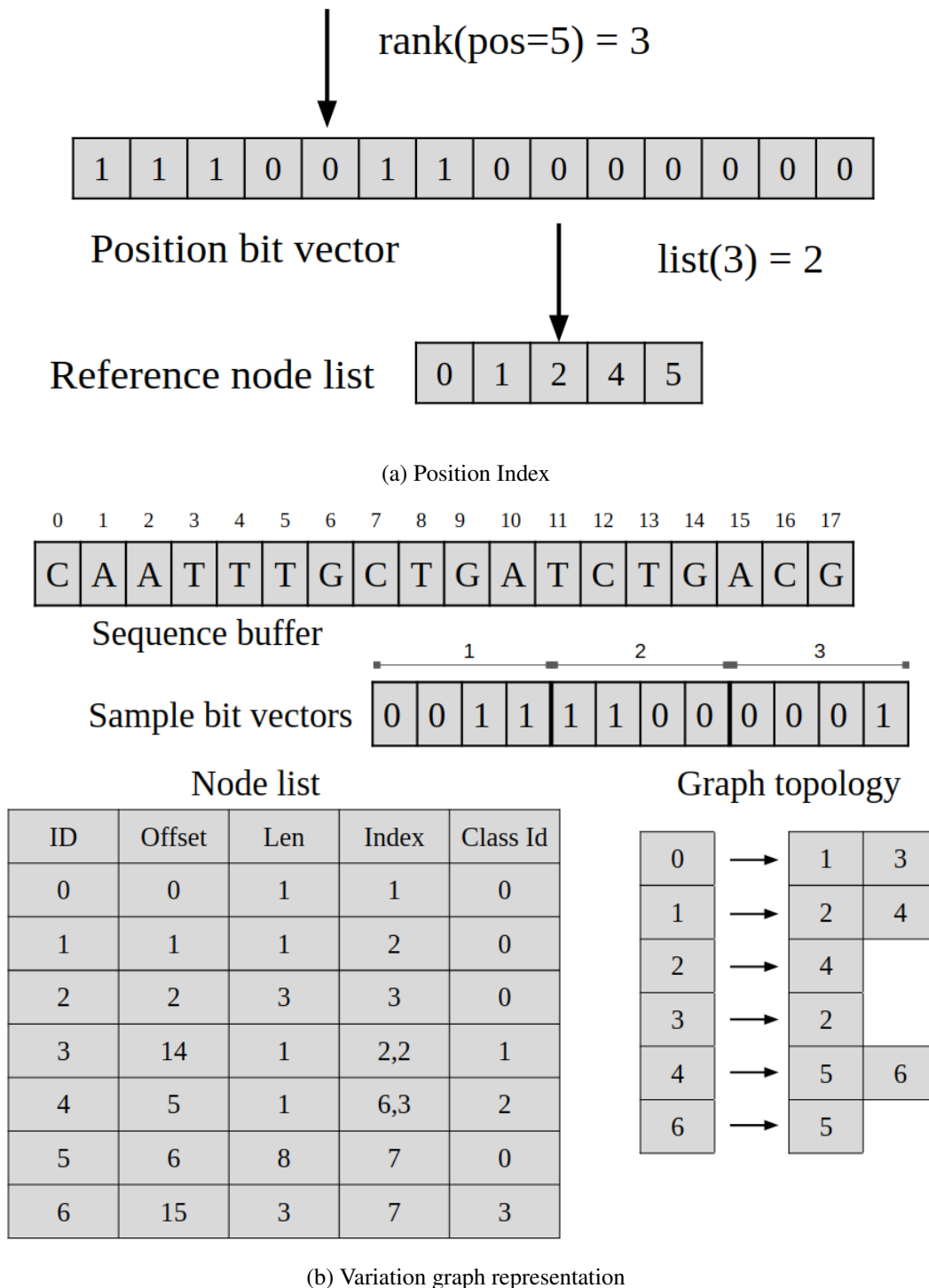


Figure 6: Position index and variation graph representation in VariantStore for the sample graph from Figure 5. (a) Shows the query operation for finding the node at position 5 in the sequence. (b) Phasing information is omitted from the node list for simplicity of the figure. In implementation, phasing information using three bits for each sample in each node.