

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

May 6, 2020

Abstract

The ability to efficiently query genomic variants from thousands of samples is critical to achieving the full potential of many medical and scientific applications such as personalized medicine. Performing variant queries based on coordinates in the reference or sample sequences is at the core of these applications. Efficiently supporting variant queries across thousands of samples is computationally challenging. Most solutions only support queries based on the reference coordinates and the ones that support queries based on coordinates across multiple samples do not scale to data containing more than a few thousand samples. We present VariantStore, a system for efficiently indexing and querying genomic variants and their sequences in either the reference or sample-specific coordinate systems. 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 0.002 – 3 seconds using memory only 10% of the size of the full representation.

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

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

24 performed over the past decade involving genomic variation, the ability to scale these studies to large-scale
25 data available today and in the near future is still limited.

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) [13].

30 A common task is to identify all the samples with a given pattern of variants or identify all samples that have
31 variants in a given gene. These tasks are translated into variant queries that require finding variants, samples
32 with variants, or sample sequences between two positions in a chromosome. Applications often compare
33 variants across multiple samples and need to perform the above queries based on sample-specific coordinate
34 systems. A coordinate system is a system that uniquely identifies the positions of variants in a given genome.
35 Each sample in variation data has a coordinate system that can be different from the reference coordinate
36 system. Variants can appear at different positions in a sample coordinate system compared to the reference
37 due to insertions and deletions (indels).

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

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

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

49 Supporting variant queries based on multiple sample coordinate systems requires maintaining a function per
50 sample that can map a position in the reference coordinate (as present in VCF files) to the sample coordinate.

51 Maintaining thousands of such functions requires storing and accessing an order of magnitude more data
52 than only indexing variants based on a single reference coordinate system. Efficiently supporting thousands
53 of coordinate systems adversely affects the memory footprint and computational complexity of the system
54 making this problem much more challenging. This limits the scalability of variant indexes that support
55 multiple coordinate systems to variation data containing only a few thousands samples.

56 VG toolkit [14] is one of the most widely used tools to represent genomic variation data and it also sup-
57 ports multiple coordinate systems. It encodes genomic variants from multiple samples in a graph, called a
58 variation graph. A variation graph is a sequence graph where each node represents a sequence and a set of
59 nodes through the graph, known as a path, embeds the complete sequence corresponding to the reference or
60 a sample. Each node on a path is assigned a position indicating the location of the sequence in the coordi-
61 nate system of the path. A node can be assigned multiple positions based on the number of paths that pass
62 through the node. The variation graph enables read alignment against multiple sample sequences contain-
63 ing variants simultaneously and avoids mapping biases that arise when mapping reads to a single reference
64 sequence [14–18].

65 VG toolkit stores each sample path as a list of nodes in the graph and maintains a separate index corre-
66 sponding to the coordinates of the reference and samples. Storing a separate list of nodes for each sequence
67 impedes the scalability of the representation for storing variation from thousands of samples. Moreover,
68 variants are often shared among samples, so storing a list of nodes for each sample path introduces redun-
69 dancy in the representation. VG toolkit is designed to optimize read alignment and uses sequence-based
70 indexes for alignment. It can not be directly used for variant queries that require an index based on the
71 position of variants in multiple sequence coordinates. Finally, the VG toolkit representation does not store
72 phasing information contained in VCF files, which is required in many analyses.

73 Multiple solutions have been proposed that efficiently index variants and support a subset of the variant
74 queries described above. GQT [19] was the first tool that proposed a sample-centric index for storing and
75 querying variants. It stores variants in compressed bitmap indexes and supports efficient comparisons of
76 sample genotypes across many variant positions. BGT [20] and GTC [21] proposed variant-centric indexes
77 that store variants in compressed bit matrices. They support queries for variants in a given region and allow
78 filtering returned variants based on subsets of samples. The SeqArray library [22] is another variant-centric

79 tool for the R programming language to store and query variants.

80 However, these tools index variants only in the reference coordinate system and do not support variant
81 queries in a sample coordinate system. Supporting multiple coordinate systems is a much harder problem to
82 tackle. Furthermore, these tools do not store the reference sequence and cannot be directly used to query and
83 compare sample sequences in a given region. Other tools have been proposed that use traditional database
84 solutions, such as SQL and NoSQL [23–25]. However, they have proven prohibitively slow to index and
85 query collections of variants.

86 We present VariantStore, a system for efficiently indexing and querying genomic information (genomic
87 variants and phasing information) from thousands of samples containing millions of variants. VariantStore
88 supports querying variants occurring between two positions across a chromosome based on the reference or
89 a sample coordinate. VariantStore bridges the gap between the tools that are space-efficient and fast but only
90 support reference-based queries (e.g., GTC [21]) and VG, which maintains multiple coordinate systems by
91 storing variants in a variation graph but fails to scale to thousands of samples. We show this by indexing
92 variants from both the 1000 genomes and TCGA projects ($> 8K$ samples), and show that VariantStore
93 is faster than VG and takes less memory and disk space. VariantStore performs variant queries based on
94 sample coordinates in less than a second. Furthermore, we have designed VariantStore to efficiently scale
95 out of RAM to storage devices in order to cater to the ever increasing sizes of available variation data by
96 performing memory-efficient construction and query.

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

104 To perform index construction and query efficiently in terms of memory, we partition the variation graph
105 into small chunks (usually a few MBs in size) based on the reference coordinates and store variation graph
106 nodes in these chunks. During construction, an active chunk is always maintained in memory in which new

107 nodes are added, and once it reaches its capacity we compress and serialize it to disk and create a new active
108 chunk. The nodes in and across chunks are ordered based on the reference coordinate since they are created
109 based on variants in the VCF file which are themselves ordered by the reference coordinate. During a query,
110 we only load the chunks in memory that contain the nodes corresponding to the query range.

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

118 **Results**

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

121 To calibrate our performance numbers, we compare VariantStore against VG toolkit [26]. VG toolkit rep-
122 resents variants in a variation graph and supports multiple coordinate systems but does not support variant
123 queries. Therefore, we only compare the construction performance and disk space against VG toolkit.

124 **Data.** We use 1000 Genomes Phase 3 data [27] and three of the biggest projects from TCGA in terms of the
125 number of samples, Ovarian Cancer (OV), Lung Adenocarcinoma (LUAD), and Breast Invasive Carcinoma
126 (BRCA), for our evaluation. 1000 Genomes data contains more variants compared to the TCGA data but
127 TCGA data contains more samples. 1000 Genomes data contains a separate VCF file for each chromosome
128 containing variants from thousands of samples. The number of samples in each file is $\approx 2.5K$. The variants
129 in 1000 Genomes project are based on GRCh37 reference genome. The TCGA data contains a separate
130 VCF file for each sample. The OV, LUAD, and BRCA projects contain 2436, 2680, and 4319 VCF files
131 containing both normal and tumor variants respectively. For each project in TCGA, we first merged VCF
132 files using the BCF tool merge command [28] and created a separate VCF file for each chromosome. The
133 variants in the TCGA project are based on GRCh38 reference genome.

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 variant index on the 1000 Genomes and TCGA (OV, LUAD, and BRCA) data using VariantStore and VG toolkit. *Space for the XG index that does not contain any path information. 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.

134 **Index construction.** The total time taken to construct the variation graph representation and index includes
135 the time taken to read and parse variants from the VCF file, construct the variation graph representation and
136 indexes, and serialize the final representation to disk. For VariantStore, the reported time includes the time
137 to create and serialize the position index. The space reported for VariantStore is the sum of the space of the
138 variation graph representation and position index.

139 For VG toolkit, creating a variation graph representation with multiple coordinate systems (or sample path
140 annotations) requires creating two indexes, XG and GBWT index. The XG index is a succinct representa-
141 tion of the variation graph without path annotations that allows memory- and time-efficient random access
142 operations on large graphs. The GBWT (graph BWT) is a substring index for storing sample paths in the
143 variation graph. We first create the variation graph representation using the “construct” command including
144 all sample path annotations. We then create the XG index and GBWT index from the variation graph rep-
145 resentation to create an index with all sample path annotations in the variation graph. For VG toolkit, the
146 reported time includes the time to create and serialize the XG and GBWT indexes. The space reported for
147 VG toolkit is the sum of the space of the XG and GBWT indexes. VG toolkit could not build GBWT index
148 on TCGA data even after running for more than a day. We only report space for the XG index (which does
149 not contain any sample path annotations) for TCGA data.

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

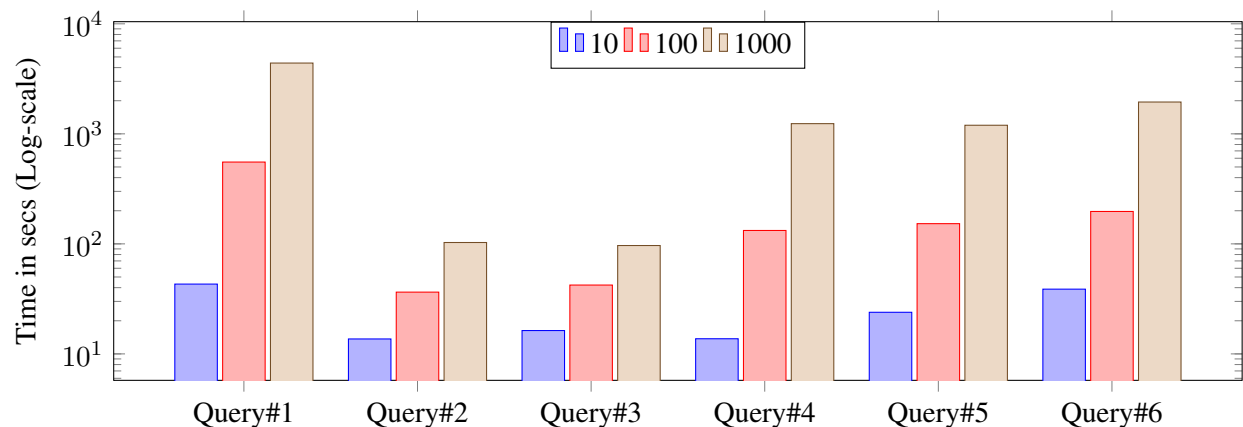
155 The performance of VariantStore and VG toolkit for constructing the index on the 1000 Genomes and TCGA
156 data is shown in Table 1.

157 VariantStore is $3\times$ faster, takes 25% less disk space, and $3\times$ less peak RAM than VG toolkit. For the
158 TCGA data, VG toolkit could not build GBWT index embedding all sample paths. However, even the space
159 needed for the XG index (not embedding sample paths) is $\approx 3.3\times$ larger than the VariantStore representation
160 containing all sample paths.

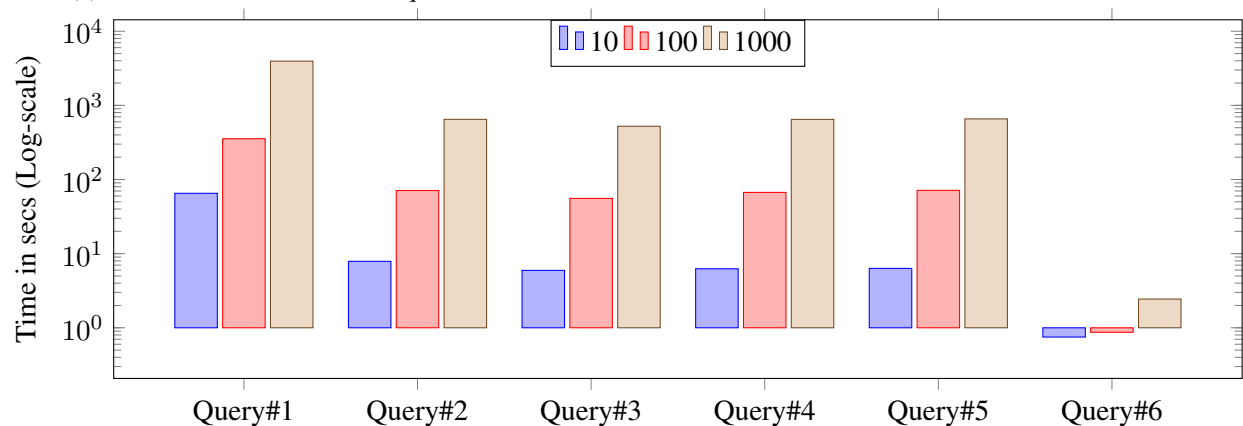
161 **Query throughput.** We measured the query throughput for all six queries mentioned above. To show the
162 robustness of query efficiency on different data we evaluate on two different chromosome indexes from two
163 different projects. We chose Chromosome 2 which is one of the bigger chromosomes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

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

174 To measure the query throughput, we only load the position index and graph topology in memory and the
175 variation graph representation is kept on disk before performing the query. Keeping the variation graph
176 representation on disk keeps the peak memory usage low and all disk accesses are performed during the
177 query to load appropriate node chunks.



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

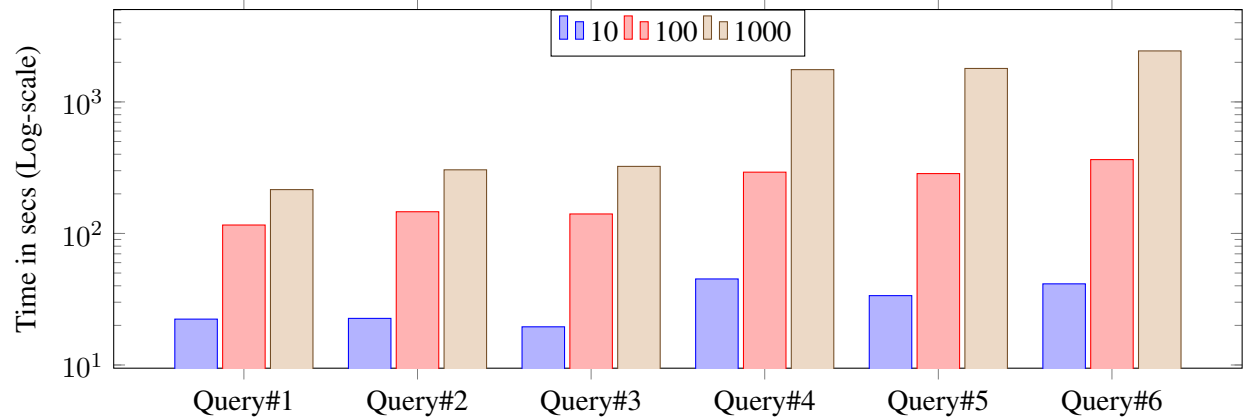


(b) Time for 10, 100, and 1000 queries on Chromosome 22 index in VariantStore for TCGA LUAD data.

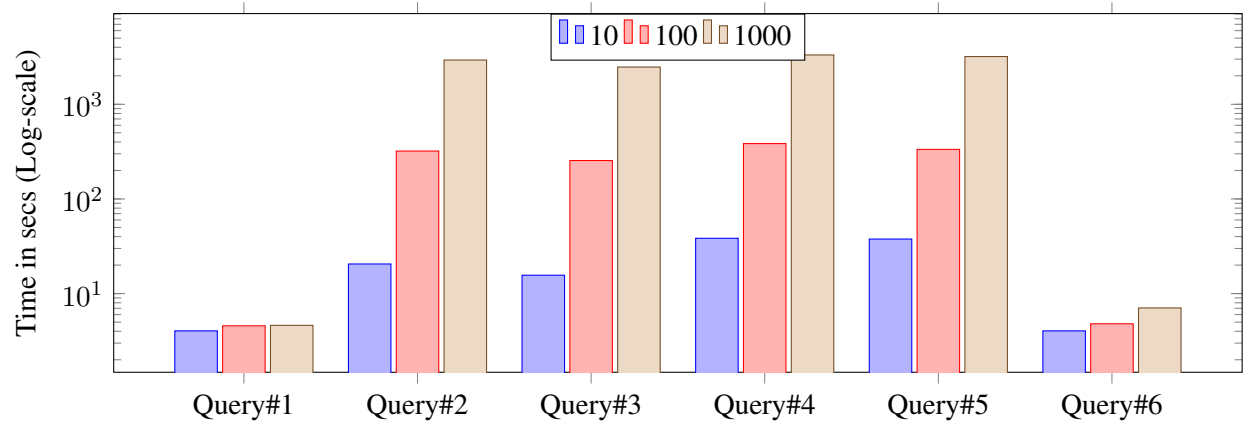
Figure 1: Time reported is the total time taken to execute 10, 100, and 1000 queries. For all queries the query length is fixed to $\approx 42K$. Query#1: Closest variant, Query#2: Seq in ref coordinate, Query#3: Seq in Sample coordinate, Query#4: Sample variants in ref coordinate, Query#5: Sample variants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

178 The query throughput on 1000 Genomes data is shown in Figures 1a and 2a. For all query types, the
 179 aggregate time taken to execute queries increases linearly with the number of queries. Finding the sequence
 180 corresponding to a sample in a region (Queries #2, #3) takes less time compared to finding variants in a
 181 region (Queries #4, #5, #6). Finding the sequence takes less time because it involves traversing the sequence
 182 specific path in the region and reconstructing the sequence. However, finding variants in a given region takes
 183 more time because it involves an exhaustive search of neighbors at each node in the region to determine all
 184 the variants that are contained by a given sample or all samples.

185 The query throughput on TCGA data is shown in Figures 1b and 2b. The TCGA data has twice as many
 186 samples compared to 1000 Genomes data but there are fewer variants. This makes the variation graph much



(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 variants in ref coordinate, Query#5: Sample variants in Sample coordinate, Query#6: All variants in ref coordinate

187 sparser and queries in the position index become more expensive compared to traversing the graph between
 188 two positions. Finding all variants is the fastest query because the variation graph is very sparse and most
 189 position ranges are empty. Traversing the graph to find the sequence or variants for a sample takes similar
 190 amount of time. Finding the closest variant from a position takes the most amount of time because it involves
 191 performing multiple position-index queries to determine the closest variant.

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

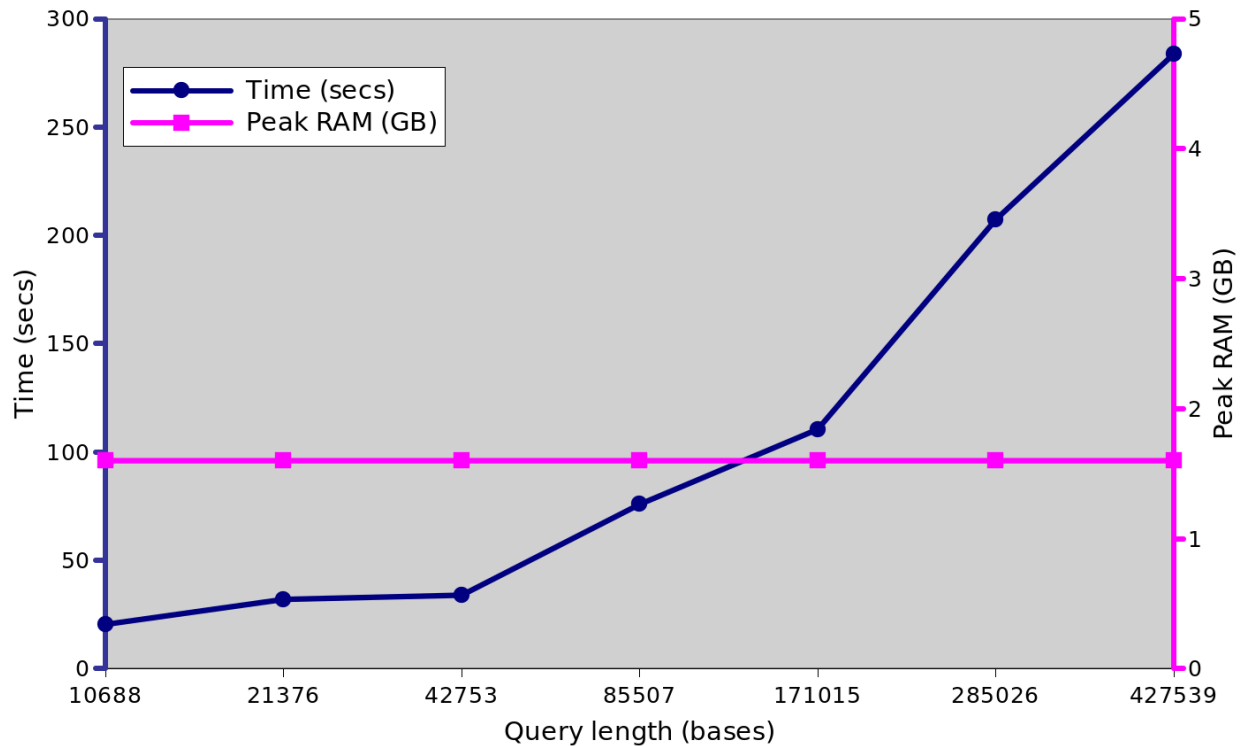


Figure 3: Time (seconds) and peak memory usage (GB) with increasing query length for “Sample sequence in the reference coordinate” (query #2) for 1000 Genomes chromosome 22 index. The time taken increases as the query length increases. But the memory usage remains constant regardless of the query length.

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

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

205 **The effect of number of variants on query time.** We also evaluate how the number of variants in the
206 position range affects the query time. For this benchmark, we performed the “All variants in the reference
207 coordinate” query (query #6) because this query involves performing a breadth-first search in the graph to

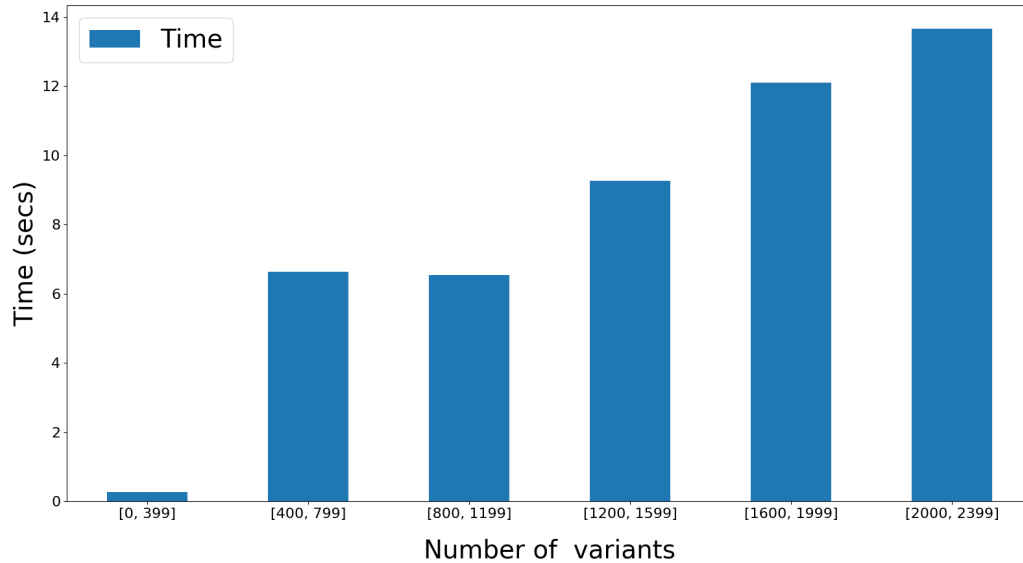


Figure 4: Time (seconds) and the number of variants in the query region for “All variants in the reference coordinate” (query #6) for 1000 Genomes chromosome 22 index. Query times are binned based on the number of variants in the query region and mean time is reported in each bin. The mean time increases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

208 determine all variants in a region and the query time depends on the number of variants in the region. To
209 perform queries on regions with different number of variants, we chose 1000 regions with a fixed size of the
210 position range ($\approx 42\text{K}$ bases) and start position chosen uniformly randomly across the chromosome.

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

214 **Comparison with a reference-based variant index.** To understand the overhead of maintaining thousands
215 of coordinate systems on the performance of VariantStore, we compare VariantStore to a variant index that
216 only indexes variants based on the reference coordinate system and supports a subset of variant queries. We
217 use GTC [29] in our evaluation as it is the fastest and smallest reference-based variant index. We compare
218 the query performance for two variant queries (Query#4 “Sample variants in ref coordinate” and Query#6
219 “All variants in ref coordinate”) supported by GTC.

220 GTC took $6\times$ less time and an order of magnitude less space to construct and store the variant index com-
221 pared to VariantStore. Furthermore, variant queries were also about an order of magnitude faster in GTC

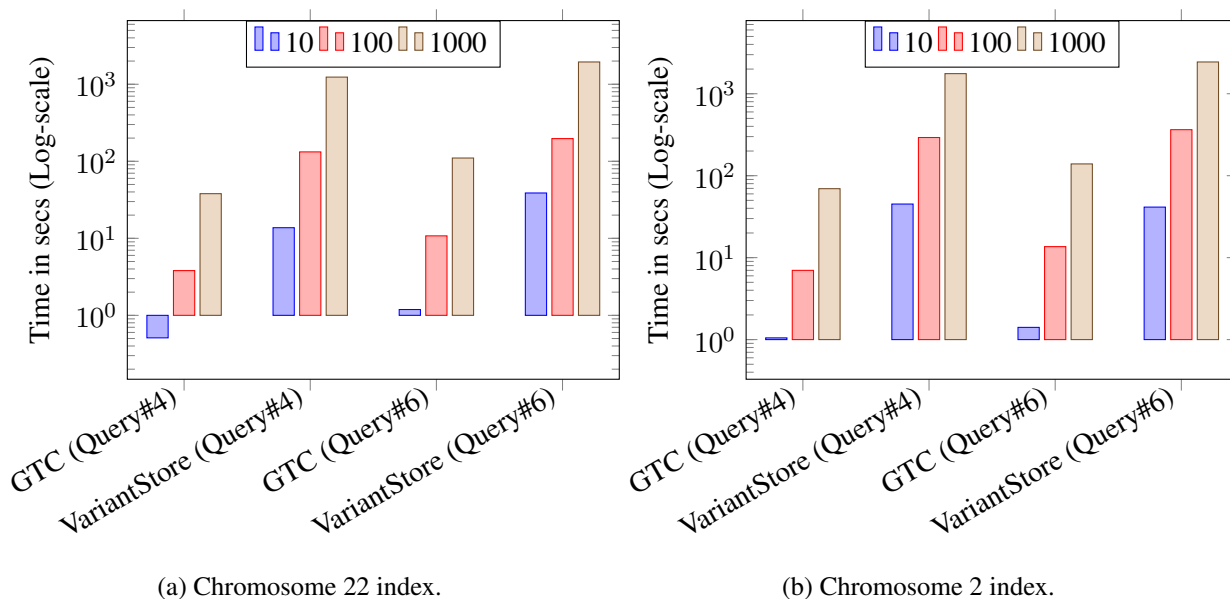


Figure 5: Time reported is the total time taken to execute 10, 100, and 1000 queries on 1000 genomes data. For all queries the query length is fixed to $\approx 42K$. Query#4: Sample variants in ref coordinate and Query#6: All variants in ref coordinate

222 (see Figure 5). The slow performance of VariantStore is because of the overhead of maintaining the variation
223 graph for representing multiple coordinate systems. During index creation, adding a variant requires split-
224 ting a reference node, adding the variant node in the graph, and updating the mapping function for all the
225 samples corresponding to the variant. To query, we first map the variant position to a node in the graph using
226 the position-index and then traverse the path in the graph to answer queries. On the other hand, adding and
227 querying variants can be performed fairly efficiently using compressed bit vectors in reference-only variant
228 indexes. If an application only requires querying variants based on a single reference coordinate system
229 then GTC offers a space-efficient and faster alternative, but it does not support queries using per-sample
230 coordinate systems or general genome graph traversals.

231 **Experimental hardware.** All experiments on 1000 Genomes data were performed on an Intel Xeon CPU
232 E5-2699A v4 @ 2.40GHz (44 cores and 56MB L3 cache) with 1TB RAM and a 7.3TB HGST HDN728080AL
233 HDD running Ubuntu 18.04.2 LTS (Linux kernel 4.15.0-46-generic) and were run using a single thread.
234 Benchmarks for TCGA data were performed on a cluster machine running AMD Opteron Processor 6220
235 @ 3GHz with 6MB L3 cache. TCGA data was stored on a remote disk and accessed via NFS.

236 **Discussion**

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

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

247 The variation graph representation in VariantStore is smaller and more efficient to construct than the rep-
248 resentation in VG toolkit. It can be further used in read alignment as a replacement to the variation graph
249 representation in the VG toolkit.

250 While the current implementation does not support adding new variants or updating the reference sequence
251 in an existing VariantStore, this is not a fundamental limitation of the design. In future, we plan to extend
252 the immutable version of VariantStore to support dynamic updates following the LSM-tree design [30].

253 **Methods**

254 **Variation graph**

255 A variation graph (VG) [14] (also defined as a genome graph in Kim et al. [17] and Rakocevic et al. [18]) is
256 a directed, acyclic graph (DAG) $G = (N, E, P)$ that embeds a set of DNA sequences. It comprises of a set
257 of nodes N , a set of direct edges E , and a set of paths P . For DNA sequences, we use the alphabet $\{A, C,$
258 $G, T, N\}$. Each $n_i \in N$ represents a sequence $seq(n_i)$. Edges in the graph connect nodes that are followed
259 on a path. Each sample in the VCF file follows a path through the variation graph. The embedded sequence
260 given by the path is the sample sequence. Given that a variation graph is a directed graph, edges can be

261 traversed in only one direction. Although, not applicable to VariantStore, an edge can also be traversed in
262 the reverse direction when the variation graph is used for read alignment [14, 17, 18].

263 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
264 a source node n_s and a destination node n_d . Nodes on a path are assigned positions based on the coordinate
265 systems of sequences they represent [31]. The position of a node on a path is the sum of the lengths of the
266 sequences represented by nodes traversed to get to the node on the path. For a path $p = n_1 \dots n_p$, position
267 $P(n_p)$ is $\sum_{i=1}^{p-1} |seq(n_i)|$. A node in the graph can appear in multiple paths and therefore can have multiple
268 positions based on different coordinate systems.

269 An initial variation graph is constructed with a single node and no edges using a linear reference sequence
270 and a reference genome coordinate system. Variants are added to the variation graph from one or more VCF
271 files [13]. A variant is encoded by a node in the variation graph that represents the variant sequence and
272 is connected to nodes representing the reference sequence via directed edges. Each variant in the VCF file
273 splits an existing reference sequence node into two (or three in some cases) and joins them via an alternative
274 path corresponding to the variant. For example, a substitution or deletion can cause an existing reference
275 node to be split into three parts (Figure 6). An insertion can cause the reference node to be split into two
276 parts (Figure 6).

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

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

282 Much work has been done devising efficient approaches to handle multiple coordinate systems. Rand et al.
283 [31] introduced the offset-based coordinate system. VG toolkit [26] implemented the multiple coordinate
284 system by explicitly storing a list of node identifiers and node offsets for each path in the variation graph.
285 They store the list of node identifiers as integer vectors using the succinct data structure library (SDSL [32]).
286 We call this an *explicit-path representation*.

287 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.

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

291 VariantStore

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

294 The variation graph representation is divided into three components:

- 295 1. Variation graph topology
- 296 2. Sequence buffer
- 297 3. List of variation graph nodes

298 **Variation graph topology.** A variation graph constructed by inserting variants from VCF files often shows
299 high sparsity (the number of edges is close to the number of nodes). For example, the ratio of the number
300 of edges to nodes in the variation graph on 1000 Genomes data [27] is close to 1. Given the sparsity of the
301 graph, we store the topology of the variation graph in a representation optimized for sparse graphs.

302 Our graph representation uses the counting quotient filter (CQF) [33] as the underlying container for storing
303 nodes and their outgoing neighbors. The CQF is a compact representation of a multiset S . Thus a CQF
304 supports inserts and queries. A query for an item x returns the number of instances of x in S . The CQF
305 uses a variable-length encoding scheme to count the multiplicity of items. In our graph representation, we
306 encode the node as the item and id of the outgoing neighbor as the count of the item in the CQF.

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

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

311 We use the version of the counting quotient filter with no false-positives to map a node to its outgoing
312 edge(s). We store the node id as the key in the CQF and if there is only one outgoing edge we encode the
313 outgoing neighbor id as the count of the key. If there are more than one outgoing edges we use indirection.
314 We maintain a list of vectors where each vector contains a list of outgoing neighbor identifiers corresponding
315 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
316 the list of outgoing neighbor ids) as the count of the key.

317 **Sequence buffer.** The sequence buffer contains the reference sequence and all variant sequences corre-
318 sponding to each substitution and insertion variant. All sequences are encoded using 3-bit characters in an
319 integer vector from SDSL library [32, 34]. The integer vector initially only contains the reference sequence.
320 Sequences from incoming variants are appended to the integer vector. Once all variants are inserted the
321 integer vector is bit compressed before being written to disk.

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

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

327 Our representation of the list of samples is based on two observations. First, multiple samples share a variant
328 and storing a list of sample identifiers for each variant is space inefficient. Instead, we store a bit vector of
329 length equal to the number of samples and set bits corresponding to the present samples in the bit vector.
330 Second, multiple variants share the same set of samples. We define an equivalence relation \sim over the set
331 of variants. Let $E(v)$ denote the function that maps each variant to the set of samples that have the variant.
332 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
333 samples shared by variants as the *sample class*. A unique id is assigned to each sample class and nodes store
334 the sample class id instead of the whole sample class. This scheme has been employed previously by other
335 colored de Bruijn graph representation tools [35–38] for efficiently maintaining a mapping from k -mers (a
336 k -length substring sequence) to the set of samples where k -mers appear.

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

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

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

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

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

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

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

366 We create the position index using a bit vector called the position-bv of length equal to the reference se-
367 quence length and a list of node identifiers on the reference path in the increasing order by their reference
368 sequence positions. For every node in the list we set the bit corresponding to the node's position in the
369 position-bv. There is a one-to-one correspondence between every set bit in the position-bv and node posi-
370 tions in the list. We store the position-bv using a bit vector and node list as an integer vector from the SDSL
371 library [32, 34].

372 **Variation graph construction**

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

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

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

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

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

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

Algorithm 1 Pseudocode to fix sample positions in the variation graph. A node corresponding to a variant contains the list of sample identifiers that have the variant and their respective positions in sample paths. A node corresponding to the reference sequence contains the position in the reference path and optionally a list of sample identifiers if it also represents a delete variant.

```
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  $neighbor.pos[sample] = 0$  then
7:          $neighbor.pos[sample] \leftarrow node.pos[ref] + node.len + \delta[sample]$ 
8:       else
9:          $\delta[sample] \leftarrow neighbor.pos[sample] - (node.pos[ref] + node.len)$ 
10:  else
11:    for  $samples$  in  $node.samples$  do
12:       $\delta[sample] \leftarrow node.pos[sample] + node.len - node.neighbor.pos[ref]$ 
```

398 Position index construction

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

405 Variant queries

406 A query is performed in two steps. We first perform a predecessor search (largest item smaller than or equal
407 to the queried item) using the queried position in the position index to locate the node n_p with the highest

408 position smaller than or equal to the queried position pos . The predecessor search is implemented using
409 the rank operation on position-bv. For bit vector $B[0, \dots, n]$, $RANK(j)$ returns the number of 1s in prefix
410 $B[0, \dots, j]$ of B . An RRR compressed position-bv supports rank operation in constant time [32, 39]. The
411 rank of pos in position-bv corresponds to the index of the node id in the node list. Figure 7a shows a sample
412 query in the position index.

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

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

419 **Memory-efficient construction and query**

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

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

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

434 Variant queries involve traversing a path in the graph between a start and an end position or exploring the
435 graph locally around a start position. All these queries require bounded exploration of the graph for which
436 we only need to look into one or a few chunks.

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

443 **References**

- 444 [1] 1000 Genomes Project Consortium. A map of human genome variation from population-scale se-
445 quencing. *Nature*, 467(7319):1061, 2010.
- 446 [2] 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human
447 genomes. *Nature*, 491(7422):56, 2012.
- 448 [3] 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*, 526
449 (7571):68, 2015.
- 450 [4] John Lonsdale, Jeffrey Thomas, Mike Salvatore, Rebecca Phillips, Edmund Lo, Saboor Shad, Richard
451 Hasz, Gary Walters, Fernando Garcia, and Nancy Young. The genotype-tissue expression (GTEx)
452 project. *Nature Genetics*, 45(6):580, 2013.
- 453 [5] TCGA: the cancer genome atlas program, 2019. URL [https://www.cancer.gov/about-nci/organization/
454 ccg/research/structural-genomics/tcga](https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga). [Online accessed August 2019].
- 455 [6] Ananyo Choudhury, Michèle Ramsay, Scott Hazelhurst, Shaun Aron, Soraya Bardien, Gerrit Botha,
456 Emile R Chimusa, Alan Christoffels, Junaid Gamiieldien, and Mahjoubeh Sefid-Dashti. Whole-genome
457 sequencing for an enhanced understanding of genetic variation among South Africans. *Nature Com-
458 munications*, 8(1):2062, 2017.

- 459 [7] Sebastian Roskosch, Hákon Jónsson, Eythór Björnsson, Doruk Beyter, Hannes P Eggertsson, Patrick
460 Sulem, Kári Stefánsson, Bjarni V Halldórsson, and Birte Kehr. PopDel identifies medium-size dele-
461 tions jointly in tens of thousands of genomes. *bioRxiv*, page 740225, 2019.
- 462 [8] Cristian Groza, Tony Kwan, Nicole Soranzo, Tomi Pastinen, and Guillaume Bourque. Personalized
463 and graph genomes reveal missing signal in epigenomic data. *bioRxiv*, page 457101, 2019.
- 464 [9] Frank W Albert and Leonid Kruglyak. The role of regulatory variation in complex traits and disease.
465 *Nature Reviews Genetics*, 16(4):197, 2015.
- 466 [10] Karen Eilbeck, Aaron Quinlan, and Mark Yandell. Settling the score: variant prioritization and
467 mendelian disease. *Nature Reviews Genetics*, 18(10):599, 2017.
- 468 [11] Claudia MB Carvalho and James R Lupski. Mechanisms underlying structural variant formation in
469 genomic disorders. *Nature Reviews Genetics*, 17(4):224, 2016.
- 470 [12] Jerome Kelleher, Yan Wong, Anthony W Wohns, Chaimaa Fadil, Patrick K Albers, and Gil McVean.
471 Inferring whole-genome histories in large population datasets. *Nature Genetics*, 51(9):1330–1338,
472 2019.
- 473 [13] The variant call format (VCF) version 4.1 specification, 2019. URL [https://samtools.github.io/](https://samtools.github.io/hts-specs/VCFv4.1.pdf)
474 [hts-specs/VCFv4.1.pdf](https://samtools.github.io/hts-specs/VCFv4.1.pdf). [Online accessed March 2019].
- 475 [14] Erik Garrison, Jouni Sirén, Adam M Novak, Glenn Hickey, Jordan M Eizenga, Eric T Dawson, William
476 Jones, Shilpa Garg, Charles Markello, Michael F Lin, Benedict Paten, and Richard Durbin. Varia-
477 tion graph toolkit improves read mapping by representing genetic variation in the reference. *Nature*
478 *Biotechnology*, 36:875–879, 2018.
- 479 [15] Alexander Dilthey, Charles Cox, Zamin Iqbal, Matthew R Nelson, and Gil McVean. Improved genome
480 inference in the MHC using a population reference graph. *Nature Genetics*, 47(6):682, 2015.
- 481 [16] Hannes P Eggertsson, Hakon Jonsson, Snaedis Kristmundsdottir, Eiríkur Hjartarson, Birte Kehr,
482 Gisli Masson, Florian Zink, Kristjan E Hjorleifsson, Aslaug Jonasdottir, and Adalbjorg Jonasdottir.
483 GraphTyper enables population-scale genotyping using pangenome graphs. *Nature Genetics*, 49(11):
484 1654, 2017.

- 485 [17] Daehwan Kim, Joseph M Paggi, Chanhee Park, Christopher Bennett, and Steven L Salzberg. Graph-
486 based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology*,
487 37(8):907–915, 2019.
- 488 [18] Goran Rakocevic, Vladimir Semenyuk, Wan-Ping Lee, James Spencer, John Browning, Ivan J Johnson,
489 Vladan Arsenijevic, Jelena Nadj, Kaushik Ghose, and Maria C Suci. Fast and accurate genomic
490 analyses using genome graphs. *Nature Genetics*, 51:354–362, 2019.
- 491 [19] Ryan M Layer, Neil Kindlon, Konrad J Karczewski, Aaron R Quinlan, and Exome Aggregation Con-
492 sortium. Efficient genotype compression and analysis of large genetic-variation data sets. *Nature*
493 *Methods*, 13(1):63, 2016.
- 494 [20] Heng Li. BGT: efficient and flexible genotype query across many samples. *Bioinformatics*, 32(4):
495 590–592, 2015.
- 496 [21] Agnieszka Danek and Sebastian Deorowicz. GTC: how to maintain huge genotype collections in a
497 compressed form. *Bioinformatics*, 34(11):1834–1840, 2018.
- 498 [22] Xiuwen Zheng, Stephanie M Gogarten, Michael Lawrence, Adrienne Stilp, Matthew P Conomos,
499 Bruce S Weir, Cathy Laurie, and David Levine. SeqArraya storage-efficient high-performance data
500 format for WGS variant calls. *Bioinformatics*, 33(15):2251–2257, 2017.
- 501 [23] Anthony J Brookes and Peter N Robinson. Human genotype–phenotype databases: aims, challenges
502 and opportunities. *Nature Reviews Genetics*, 16(12):702, 2015.
- 503 [24] Joachim Kutzera and Patrick May. Variant-DB: A tool for efficiently exploring millions of human
504 genetic variants and their annotations. In *International Conference on Data Integration in the Life*
505 *Sciences*, pages 22–28. Springer, 2017.
- 506 [25] Geert Vandeweyer, Lut Van Laer, Bart Loeys, Tim Van den Bulcke, and R Frank Kooy. VariantDB: a
507 flexible annotation and filtering portal for next generation sequencing data. *Genome Medicine*, 6(10):
508 74, 2014.
- 509 [26] Variation graph toolkit, 2019. URL <https://github.com/vgteam/vg>. [Online accessed March 2019].

- 510 [27] IGSR: the international genome sample resource, 2019. URL <http://www.internationalgenome.org/>
511 home. [Online accessed March 2019].
- 512 [28] BCF toolkit, 2019. URL <https://samtools.github.io/bcftools/bcftools.html>. [Online accessed March
513 2019].
- 514 [29] GenoTypes compressor, 2020. URL <https://github.com/refresh-bio/GTC>. [Online accessed January
515 2020].
- 516 [30] Patrick O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth O’Neil. The log-structured merge-tree
517 (LSM-tree). *Acta Informatica*, 33(4):351–385, 1996.
- 518 [31] Knut D Rand, Ivar Grytten, Alexander J Nederbragt, Geir O Storvik, Ingrid K Glad, and Geir K
519 Sandve. Coordinates and intervals in graph-based reference genomes. *BMC Bioinformatics*, 18(1):
520 263, 2017.
- 521 [32] Simon Gog, Timo Beller, Alistair Moffat, and Matthias Petri. From theory to practice: Plug and play
522 with succinct data structures. In *13th International Symposium on Experimental Algorithms, (SEA*
523 *2014)*, pages 326–337, 2014.
- 524 [33] Prashant Pandey, Michael A Bender, Rob Johnson, and Rob Patro. A general-purpose counting filter:
525 Making every bit count. In *Proceedings of the 2017 ACM International Conference on Management*
526 *of Data*, pages 775–787. ACM, 2017.
- 527 [34] SDSL: succinct data structure library, 2019. URL <https://github.com/simongog/sdsl-sdslite>. [Online
528 accessed March 2019].
- 529 [35] Fatemeh Almodaresi, Prashant Pandey, and Rob Patro. Rainbowfish: a succinct colored de Bruijn
530 graph representation. In *17th International Workshop on Algorithms in Bioinformatics (WABI 2017)*.
531 Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- 532 [36] Prashant Pandey, Fatemeh Almodaresi, Michael A Bender, Michael Ferdman, Rob Johnson, and Rob
533 Patro. Mantis: A fast, small, and exact large-scale sequence-search index. *Cell Systems*, 7(2):201–207,
534 2018.
- 535 [37] Fatemeh Almodaresi, Prashant Pandey, Michael Ferdman, Rob Johnson, and Rob Patro. An efficient,

536 scalable and exact representation of high-dimensional color information enabled via de Bruijn graph
537 search. In *International Conference on Research in Computational Molecular Biology*, pages 1–18.
538 Springer, 2019.

539 [38] Fatemeh Almodaresi, Prashant Pandey, Michael Ferdman, Rob Johnson, and Rob Patro. An efficient,
540 scalable, and exact representation of high-dimensional color information enabled using de bruijn graph
541 search. *Journal of Computational Biology*, 27(4):485–499, April 2020. doi: 10.1089/cmb.2019.0322.
542 URL <https://doi.org/10.1089/cmb.2019.0322>.

543 [39] Rajeev Raman, Venkatesh Raman, and Srinivasa Rao Satti. Succinct indexable dictionaries with appli-
544 cations to encoding k-ary trees, prefix sums and multisets. *ACM Transactions on Algorithms (TALG)*,
545 3(4):43, 2007.

546 **Data availability**

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

548 **Acknowledgements**

549 The results published here are in whole or part based upon data generated by The Cancer Genome Atlas
550 (dbGaP accession phs000178) managed by the NCI and NHGRI. Information about TCGA can be found
551 at <http://cancergenome.nih.gov>. The 1000 Genomes data used for the analyses described in this manuscript
552 were obtained from <https://www.internationalgenome.org/>. This research is funded in part by the Gordon
553 and Betty Moore Foundation’s Data-Driven Discovery Initiative through Grant GBMF4554 to C.K. and by
554 the US National Institutes of Health (R01GM122935). We would also like to thank Shawn Baker for many
555 helpful discussions, and Guillaume Marçais and Yutong Qiu for comments on the manuscript.

556 **Competing Interests**

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

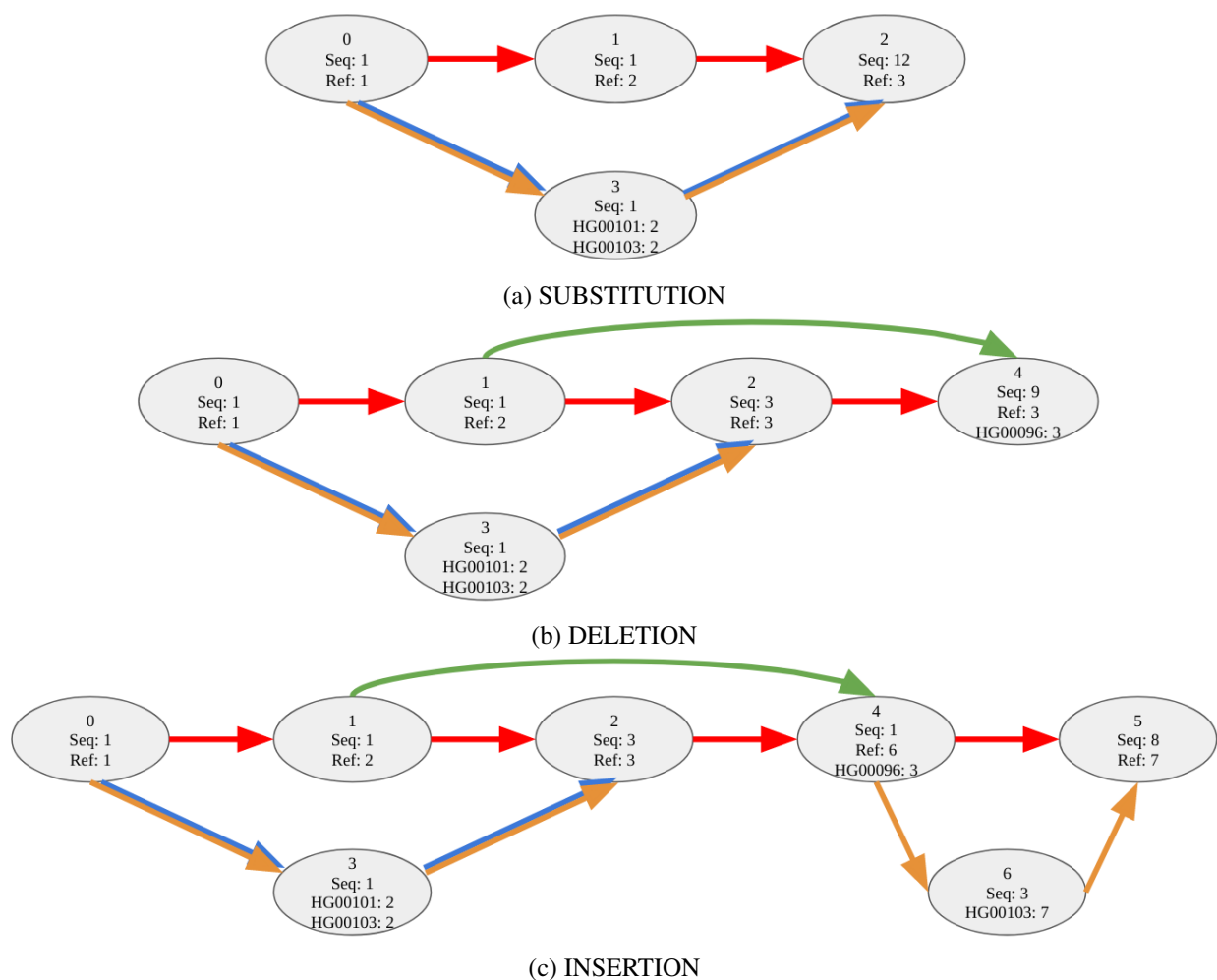
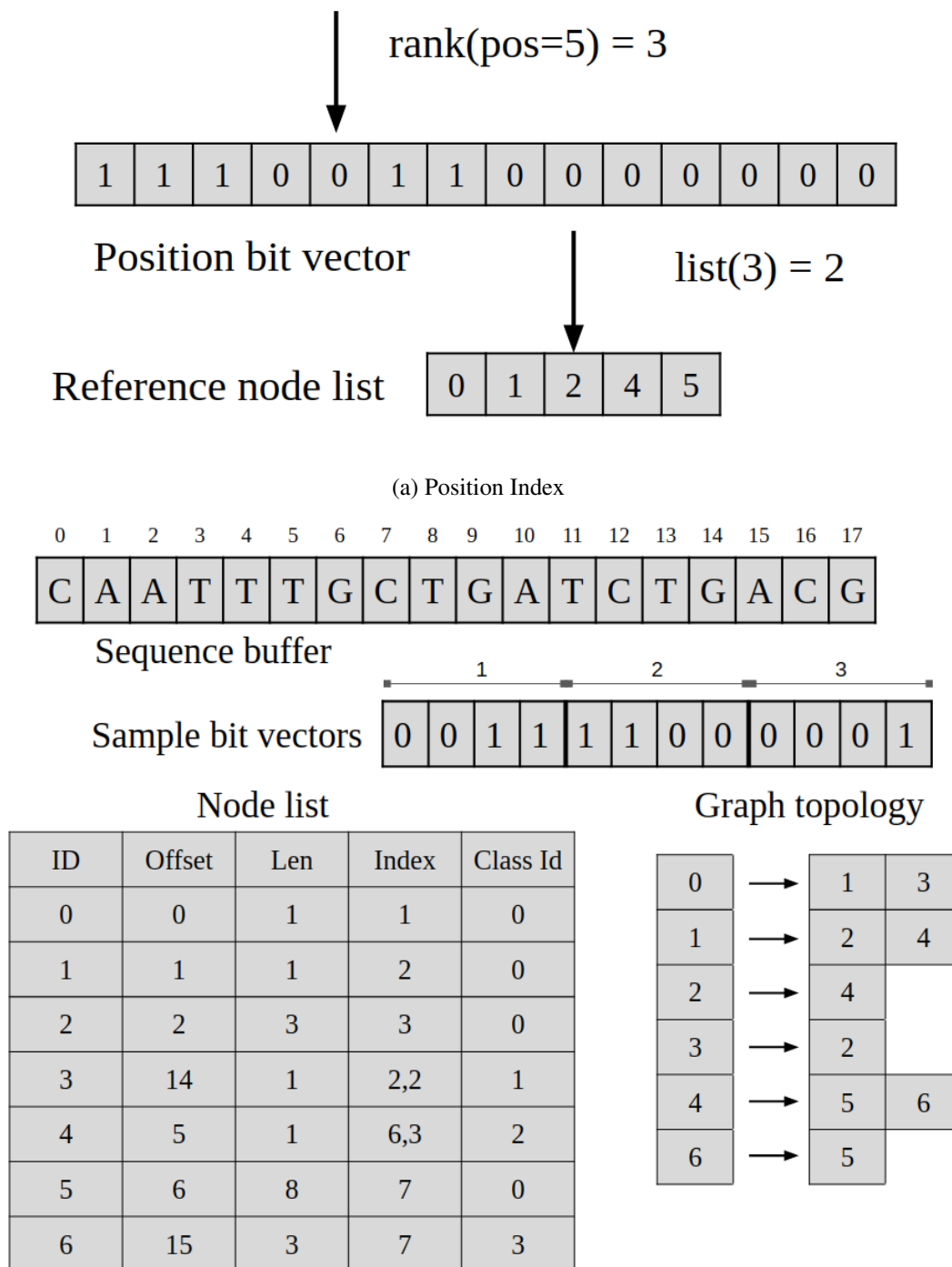


Figure 6: A variation graph with three input samples (HG00096, HG00101, HG00103) showing the encoding of substitutions, insertions, and deletions as stated in Table 2. Fig. (a) shows the substitution, (b) shows the deletion, and (c) shows the insertion. Edges are colored (or multi-colored) to show the 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.



(b) Variation graph representation

Figure 7: Position index and variation graph representation in VariantStore for the sample graph from Figure 6. (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.