

1 **Liftoff: an accurate gene annotation mapping tool**

2

3 Alaina Shumate^{1,2,*}, Steven L. Salzberg^{1,2,3,4}

4

5 1. Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD
6 21218

7 2. Center for Computational Biology, Whiting School of Engineering, Johns Hopkins
8 University, Baltimore, MD 21211

9 3. Department of Computer Science, Johns Hopkins University, Baltimore, MD 21218

10 4. Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins
11 University, Baltimore, MD 21205

12 * Correspondence: alainashumate@gmail.com

13

14 **Abstract**

15 Improvements in DNA sequencing technology and computational methods have led to a
16 substantial increase in the creation of high-quality genome assemblies of many species.
17 To understand the biology of these genomes, annotation of gene features and other
18 functional elements is essential; however for most species, only the reference genome
19 is well-annotated. One strategy to annotate new or improved genome assemblies is to
20 map or 'lift over' the genes from a previously-annotated reference genome. Here we
21 describe Liftoff, a new genome annotation lift-over tool capable of mapping genes
22 between two assemblies of the same or closely-related species. Liftoff aligns genes
23 from a reference genome to a target genome and finds the mapping that maximizes
24 sequence identity while preserving the structure of each exon, transcript, and gene. We
25 show that Liftoff can accurately map 99.9% of genes between two versions of the
26 human reference genome with an average sequence identity >99.9%. We also show
27 that Liftoff can map genes across species by successfully lifting over 98.4% of human
28 protein-coding genes to a chimpanzee genome assembly with 98.7% sequence identity.

29

30 **Availability**

31 The source code for Liftoff is available at <https://github.com/agshumate/Liftoff>

32

33 **Introduction**

34 Recent developments in DNA sequencing technology have greatly reduced the time
35 and money needed to sequence and assemble new genomes. Currently there are
36 13,420 eukaryotic genome assemblies in GenBank, of which ~10,000 have been added
37 in the last 5 years alone. The addition of new and improved genome assemblies is a
38 starting point for genetic studies of many species; however, to be maximally useful, the
39 genes and other functional elements need to be annotated. Unfortunately, the

40 annotation of new genomes has not kept pace with sequencing and assembly. This is
41 evident in GenBank, where only 3,540 of the 13,420 eukaryotic genomes have any
42 annotation at all. Eukaryotic genome annotation is a challenging, imperfect process that
43 requires a combination of computational predictions, experimental validation, and
44 manual curation. Rather than repeating this costly process for each new genome that is
45 assembled, a more scalable approach is to take the annotation from a previously-
46 annotated member of the same or closely-related species, and then map or ‘lift over’
47 gene models from the annotated genome onto the new assembly. In addition, for well-
48 studied organisms, multiple assemblies may be produced over time, and there is an
49 ongoing need to lift the annotation onto these newer, more contiguous assemblies. The
50 most well-known example of this is the human genome, but other model organisms
51 such as mouse, zebrafish (Church *et al.*, 2011), rhesus macaque (He *et al.*, 2019) ,
52 maize (Jiao *et al.*, 2017) , and many others had a series of gradually improved
53 assemblies.

54
55 Current strategies for this task use tools such as UCSC liftOver (Kuhn *et al.*, 2013) or
56 CrossMap (Zhao *et al.*, 2014) to convert the coordinates of genomic features between
57 assemblies; however, these tools only work with a limited number of species and they
58 rely only on sequence homology to find a one-to-one mapping between genomic
59 coordinates in the reference and coordinates in the target. This strategy is often
60 inadequate when converting genomic intervals, like a gene feature, rather than a single
61 coordinate. If the interval is no longer continuous in the target genome, current
62 strategies will either split the interval and map it to different locations, or map the
63 spanned interval to the target genome (Gao *et al.*, 2018). In many cases, this disrupts
64 the biological integrity of the genomic feature; for example, if the interval is split and
65 mapped to different chromosomes or strands, or spans a large genomic distance, it may
66 not be possible for it to represent a single gene feature. Furthermore, prior tools convert
67 each feature independently, so while every exon from one transcript may be lifted over
68 to a continuous interval, the combination of exons in the target genome may not
69 necessarily form a biologically meaningful transcript. Mapping each feature
70 independently also often results in multiple paralogous genes incorrectly mapping to a
71 single locus.

72
73 Here we introduce Liftoff, an accurate tool that maps annotations described in General
74 Feature Format (GFF) or General Transfer Format (GTF) between assemblies of the
75 same, or closely-related species. Unlike current coordinate lift-over tools which require a
76 pre-generated “chain” file as input, Liftoff is a standalone tool that takes two genome
77 assemblies and a reference annotation as input and outputs an annotation of the target
78 genome. Liftoff uses Minimap2 (Li, 2018) to align the gene sequences from a reference
79 genome to the target genome. Rather than aligning whole genomes, aligning only the
80 gene sequences allows genes to be lifted over even if there are many structural
81 differences between the two genomes. For each gene, Liftoff finds the alignments of the
82 exons that maximize sequence identity while preserving the transcript and gene
83 structure. If two genes incorrectly map to overlapping loci, Liftoff determines which
84 gene is most-likely mis-mapped, and attempts to re-map it. Liftoff can also find

85 additional gene copies present in the target assembly that are not annotated in the
86 reference.

87

88 Previously, we have used Liftoff to map genes to a new Ashkenazi human reference
89 genome (Shumate, Zimin *et al.*, 2020) and to an updated assembly of the bread wheat
90 genome, *Triticum aestivum* (Alonge, Shumate *et al.*, 2020). Here, in addition to
91 describing the algorithm itself, we present two more examples demonstrating the
92 accuracy and versatility of Liftoff. First, we map genes between two versions of the
93 human reference genome. Next, to demonstrate a cross-species lift over, we map
94 protein-coding genes from the human reference genome to a chimpanzee genome
95 assembly.

96

97 **Implementation**

98 Liftoff is implemented as a python command-line tool. The main goal of Liftoff is to align
99 gene features from a reference genome to a target genome and use the alignment(s) to
100 optimally convert the coordinates of each exon. An optimal mapping is one in which the
101 sequence identity is maximized while maintaining the integrity of each exon, transcript,
102 and gene. While our discussion of Liftoff here focuses on lifting over genes, transcripts,
103 and exons, it will work for any feature, or group of hierarchical features present in a GFF
104 or GTF file.

105

106 As input, Liftoff takes a reference genome sequence and a target genome sequence in
107 FASTA format, and a reference genome annotation in GFF or GTF format. The
108 reference annotation is processed with gffutils (<https://github.com/daler/gffutils>), which
109 uses a sqlite3 database to track the hierarchical relationships within groups of features
110 (e.g. gene, transcript, exon). Using pyfaidx (Shirley *et al.*, 2015), Liftoff extracts gene
111 sequences from the reference genome, and then invokes Minimap2 to align the entire
112 gene sequence including exons and introns to the target. The Minimap2 parameters are
113 set to output up to 50 secondary alignments for each sequence in SAM format. By
114 default, genes are aligned to the entire target genome, but for chromosome-scale
115 assemblies, the user can enable an option to align genes chromosome by chromosome.
116 Under that option, only those genes which fail to map to their expected chromosome are
117 then aligned to the entire genome.

118

119 In many cases, a gene has a single complete alignment to the target genome, which
120 makes finding the optimal mapping trivial. In other cases, differences between the two
121 genomes cause the gene to align in many fragmented pieces, and the optimal mapping
122 is some combination of alignments. To find this combination, Liftoff uses networkx (
123 <https://github.com/networkx/networkx>) to build a directed acyclic graph representing the
124 alignments as follows. Using Pysam (<https://github.com/pysam-developers/pysam>) to
125 parse the Minimap2 alignments, each alignment is split at every insertion and deletion in
126 order to form a group of gapless alignment blocks. Blocks not containing any part of an
127 exon are discarded, and the remaining blocks are represented by nodes in the graph.
128 Two nodes u and v are connected by an edge if the following conditions are true.

129 1) u and v are on the same chromosome or contig

130 2) u and v are on the same strand

- 131 3) u and v are in the correct 5' to 3' order
132 4) The distance from the start of u to the end of v in the target genome is no greater
133 than 2 times that in the reference genome

134

135 Each node is assigned a weight equal to the number of mismatches within exons
136 (mismatches in introns are not counted), and each edge is assigned a weight equal to
137 the length of gaps within exons spanned by that edge. A source and sink are added to
138 the graph representing the start and end of the gene respectively, and the shortest path
139 from source to sink is found using Dijkstra's algorithm (Dijkstra and Others, 1959)
140 where the weight function between two nodes u and v is

141

$$142 \quad \frac{weight_u + weight_v}{2} + weight_{edge}$$

143

144 The shortest path represents the combination of aligned blocks that is concordant with
145 the original structure of the gene and minimizes the number of mismatches and indels
146 within exons. The alignments in this path define the final placement of the gene. Using
147 the coordinates of the aligned blocks in the shortest path, the coordinates of each exon
148 are converted to their respective coordinates in the target genome. A simple example of
149 this process is shown in **Figure 1**, which illustrates lifting over a 5-exon transcript from
150 the human reference genome (GRCh38) to a chimpanzee genome (PTRv2). This gene
151 has a large intronic deletion in PTRv2 and does not have an end-to-end alignment, but
152 it can still be successfully lifted over using our algorithm.

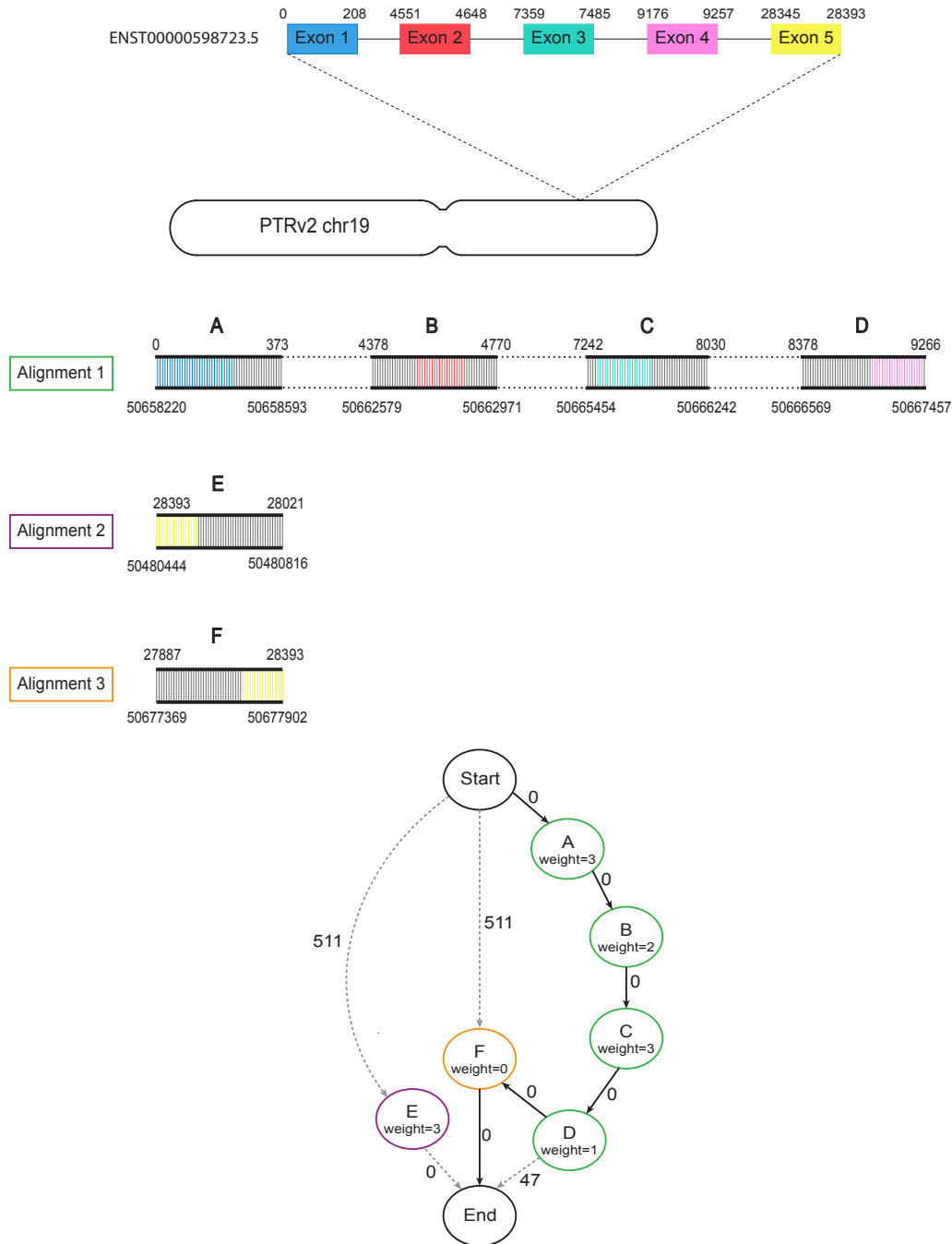


Figure 1. Example of the lift-over process. Diagram showing the steps taken by Liftoff when mapping human transcript ENST00000598723.5 to the chimpanzee (PTRv2) homolog on chromosome 19. Minimap2 produces 3 partial alignments of this gene to PTRv2. Alignment 1 (green) has 4 gapless blocks containing exons 1-4 which are represented by nodes A-D in the graph. The dashed lines in between blocks of the alignment represent gaps/introns. Alignments 2 (purple) and 3 (orange) each have 1 gapless block containing exon 5 represented by nodes E and F respectively. Node E is not on the same strand as alignments 1 and 2 and is therefore only connected to the start and end. The node weights correspond to the number of mismatches in exons and the edge weights are the number of unaligned exon bases between two nodes. The shortest path (A,B,C,D,F) is shown with bold arrows and contains complete alignments of all 5 exons with a total of 9 mismatches and 0 gaps.

153 One of the main challenges with gene annotation lift over is correctly mapping
154 homologous genes from multi-gene families. Two different genes may optimally map to
155 the same locus if they are identical or nearly identical. To handle this situation, after
156 Liftoff maps all genes to their best matches, it checks for pairs of genes on the
157 reference genome that have incorrectly mapped to overlapping (or identical) locations
158 on the target genome, and it then attempts to find another valid mapping for one of the
159 genes. Liftoff first tries to remap the gene with the lower sequence identity. If the genes
160 mapped with the same sequence identity, Liftoff considers the neighboring genes and
161 tries to remap the gene that appears out of order according to the reference annotation.
162 When remapping the gene, Liftoff rebuilds the graph of aligned blocks excluding any
163 blocks that overlap the homologous gene. The shortest path through this new graph
164 represents the best mapping for this gene that does not overlap its homolog. If another
165 valid mapping does not exist, the gene with lower identity is considered unmapped. This
166 process is repeated until there are no genes mapped to overlapping loci. Liftoff then
167 outputs a GFF file with the coordinates on the target genome of all of the features from
168 the original annotation, and a text file with the IDs of any genes that could not be lifted
169 over.

170

171 Note that differences in the genome sequences themselves may result in Liftoff
172 mapping a gene to a paralogous location. For example, consider a gene family with 5
173 members on the reference genome but only 4 members on the target. The fifth gene
174 might simply be unmapped, but if the target has a paralogous copy elsewhere, *and* if
175 that copy is not matched by a homolog on the reference, then Liftoff will map the fifth
176 gene to the paralogous location.

177

178 **Annotating Extra Gene Copies**

179 Another feature unique to Liftoff is the option to find additional copies of genes in the
180 target assembly not annotated in the reference. With this option enabled, Liftoff maps
181 the complete reference annotation first, and then repeats the lift-over process for all
182 genes. An extra gene copy is annotated if another mapping is found that does not
183 overlap any previously-annotated genes, and that meets the user-defined minimum
184 sequence identity threshold. The lift-over procedure is repeated until all valid mappings
185 have been found.

186

187 We recently used Liftoff with this feature enabled to annotate our improved assembly of
188 the bread wheat genome, which contains 15.07 gigabases of anchored sequence
189 compared to 13.84 in a previous reference genome (Alonge, Shumate *et al.*, 2020). In
190 addition to successfully mapping 100,839 of the 105,200 reference genes to this large
191 and complex genome, we found 5,799 additional gene copies using a strict sequence
192 identity threshold of 100%.

193

194 **Results**

195 Here we demonstrate Liftoff's ability to lift an annotation to an updated reference
196 genome by lifting genes from the two most recent versions of the human reference
197 genome, GRCh37 and GRCh38. We also demonstrate Liftoff's ability to lift genes
198 between genomes of closely-related species by lifting genes from GRCh38 to the

199 chimpanzee genome Clint_PTRv2. To assess the accuracy of Liftoff in each example,
200 we evaluate both the sequence identity and order of mapped genes.

201

202 **GRCh37 to GRCh38**

203 We attempted to map all protein-coding genes and lncRNAs on primary chromosomes
204 in the GENCODE v19 annotation (Harrow *et al.*, 2012) from GRCh37 to GRCh38. Out
205 of 27,459 genes, we successfully mapped 27,424 (99.87%). We consider a gene to be
206 successfully mapped if at least 50% of the reference gene maps to the target assembly.
207 An overwhelming majority of the gene sequences in GRCh38 were nearly identical to
208 the sequences in GRCh37, with an average sequence identity in exons of 99.97%
209 **(Figure 2)**.

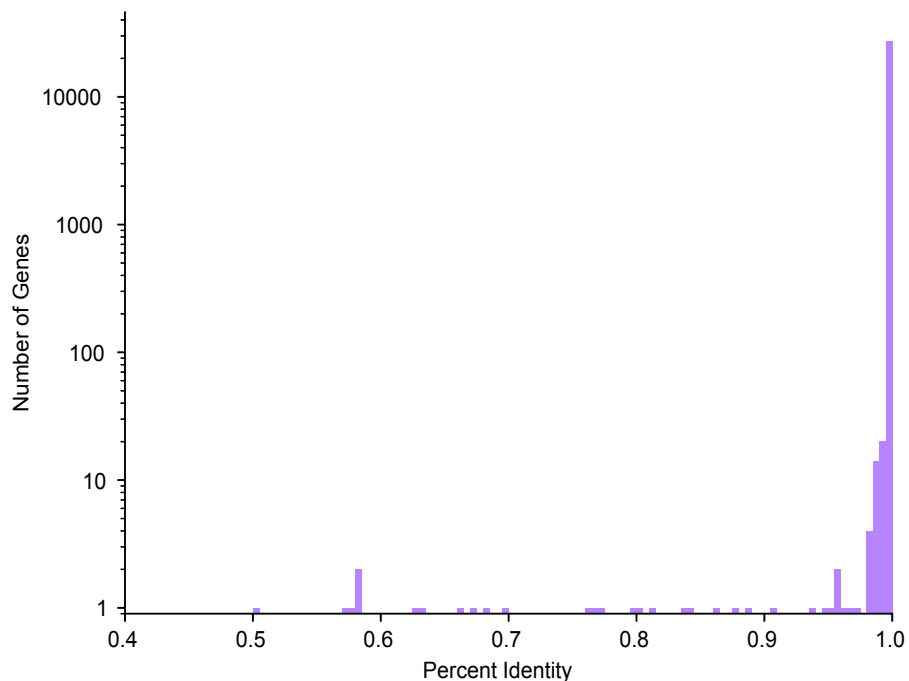


Figure 2. Distribution of GRCh37 and GRCh38 sequence identity. Histogram showing the distribution of exon sequence identity of protein-coding and lncRNA genes in GRCh37 and GRCh38. Log scale used to make the counts of just 1 or 2 genes visible; all bins below 97% identity contain at most 2 genes.

210

211 To visualize the co-linearity of the gene order between the two assemblies, we plotted
212 each gene as a single point on a 2D plot where the X coordinate is the ordinal position
213 of the gene in GRCh37 and the Y coordinate is the ordinal position in GRCh38 **(Figure**
214 **3)**.

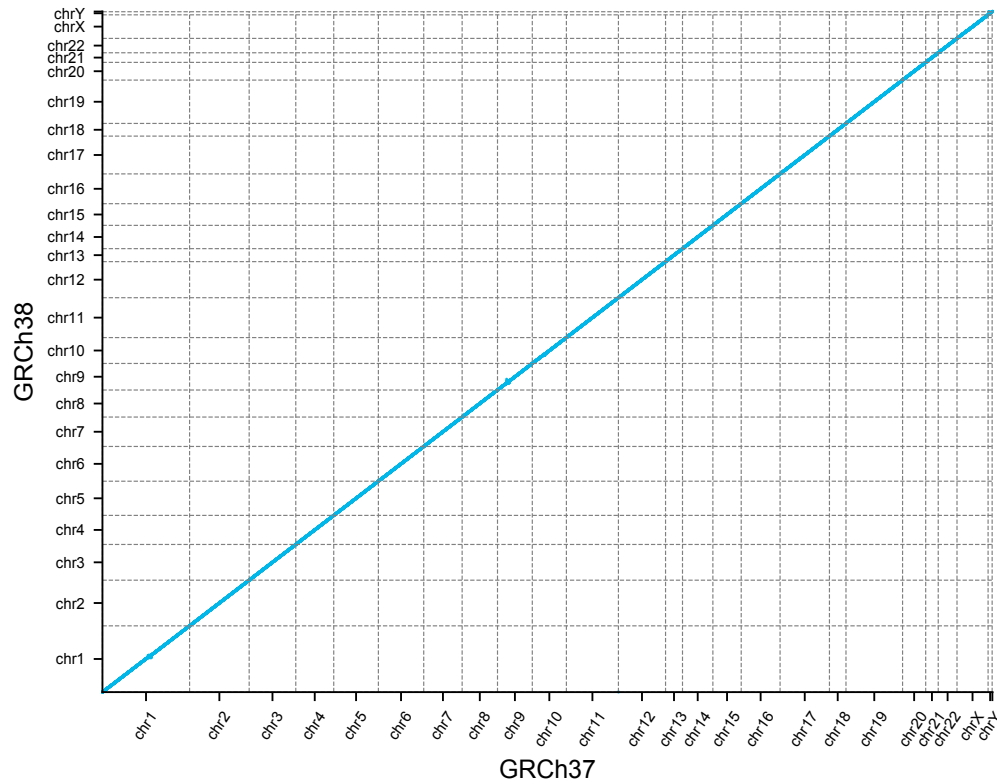


Figure 3. GRCh37 and GRCh38 gene order. Dot plot showing the ordinal position of each gene in GRCh37 on the x-axis and the ordinal position in GRCh38 on the y-axis.

215

216

217 The gene order appears perfectly co-linear; however, there are some exceptions not
218 visible at the scale of the whole genome. To calculate the number of genes out of order
219 in GRCh38 with respect to GRCh37, we first sorted the X,Y points by X, and then found
220 the length of the longest increasing subsequence in Y. The longest increasing
221 subsequence in Y represents the genes in GRCh38 that are in order with respect to
222 GRCh37, and those points not belonging to this subsequence are genes which are out
223 of order. With this process we found 305 genes (1.1%) occurring in a different relative
224 position in GRCh38 with respect to GRCh37.

225

226 **GRCh38 to PTRv2**

227 We attempted to map all protein-coding genes on chromosomes 1-22 and chromosome
228 X in the GENCODE v33 annotation (Frankish *et al.*, 2019) from GRCh38 to an
229 assembly of the chimpanzee (*Pan troglodytes*), PTRv2. Out of 19,878 genes, we were
230 able to map 19,555 (98.38%). The average sequence identity in exons of successfully
231 mapped genes was 98.70% (**Figure 4**).

232

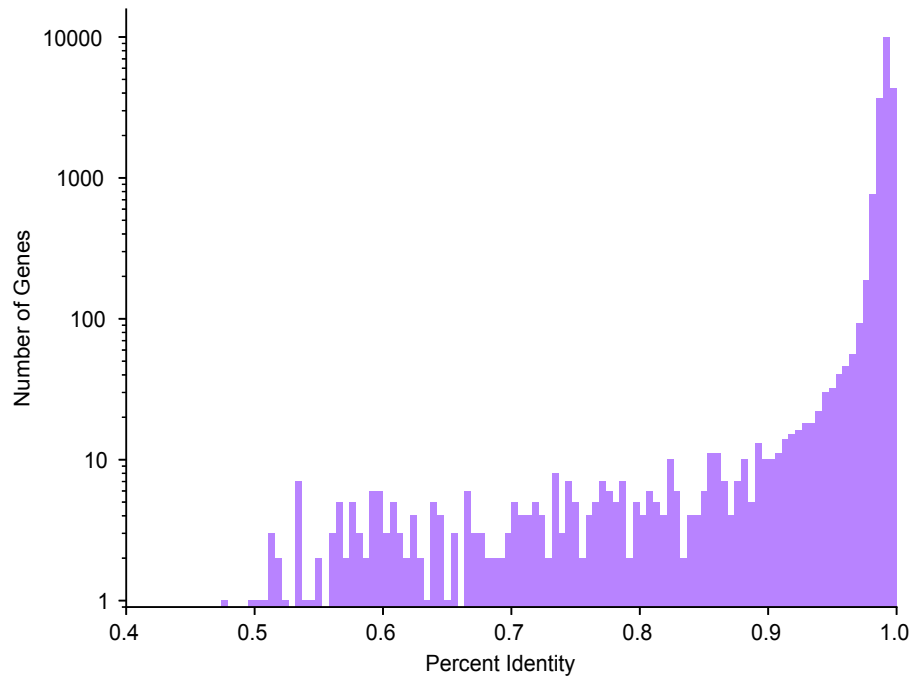


Figure 4. Distribution of GRCh38 and PTRv2 sequence identity. Histogram showing the distribution of exon sequence identity of protein-coding genes in GRCh38 and PTRv2. Note that the y-axis is shown on a log scale, as in Figure 2.

233

234

235 As was done with the GRCh37 to GRCh38 lift-over, we compared the gene order in
236 GRCh38 to that in PTRv2 and found 2,172 genes in PTRv2 to be in a different relative
237 position. Some of these ordinal differences are visible at the whole-genome scale
238 (**Figure 5**) including 4 large regions on the chimpanzee homologues of chromosomes 4,
239 5, 12, and 17 where the gene order is inverted due to large-scale chromosomal
240 inversions.

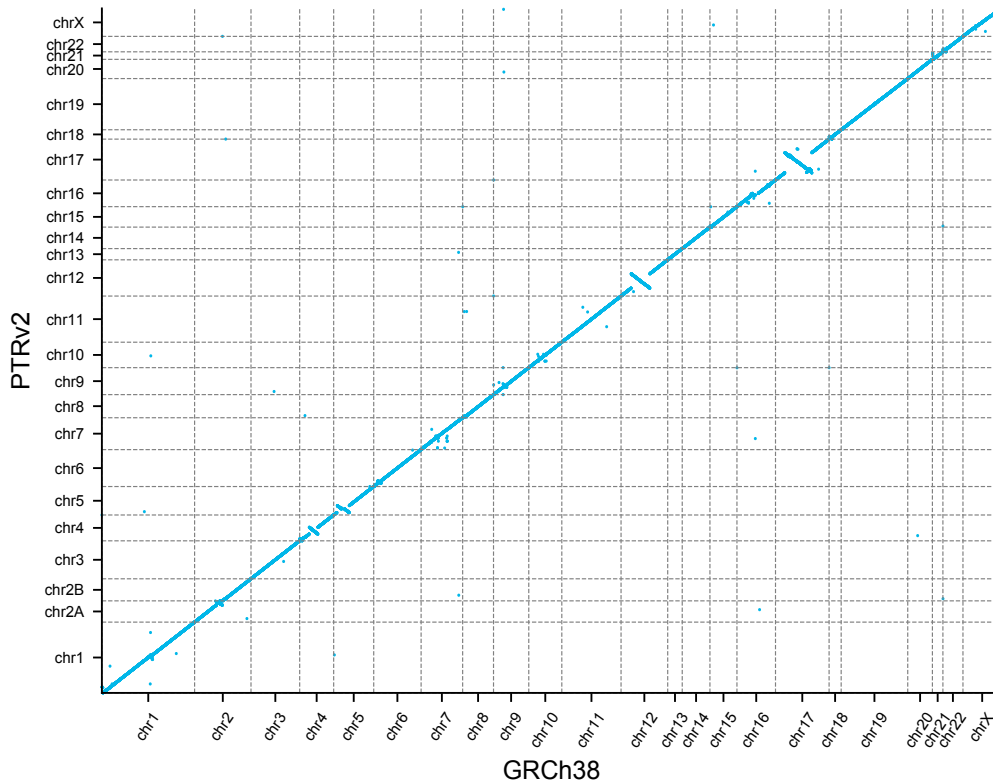


Figure 5. GRCh38 and PTRv2 gene order. Dot plot showing the ordinal position of each gene in GRCh38 on the x-axis and the ordinal position in PTRv2 on the y-axis.

241

242

243 Discussion

244 The rapidly growing number of high-quality genome assemblies has greatly increased
245 our potential to understand sequence diversity, but accurate genome annotation is
246 needed to understand the biological impact of this diversity. Rather than annotating
247 genomes *de novo*, we can take advantage of the extensive work that has gone into
248 creating reference annotations for many well-studied species. We developed Liftoff as
249 an accurate tool for transferring gene annotations between genomes of the same or
250 closely-related species. Unlike current coordinate lift-over strategies which only consider
251 sequence homology, Liftoff considers the constraints between exons of the same gene
252 and constraint that distinct genes need to map to distinct locations.

253

254 We showed that we were able to lift over nearly all genes from GRCh37 to GRCh38.
255 The gene sequences and order are very similar between the two assemblies, with an
256 average sequence identity of >99.9% and only 305 genes appearing in a different order.
257 GRCh38 fixed a number of mis-assemblies and single base errors present in GRCh37
258 (Guo *et al.*, 2017), so it is not unexpected that the gene sequence and order are not
259 entirely identical. This demonstrates Liftoff's ability to accurately annotate an updated
260 reference assembly, making it a useful tool as reference assemblies are continuously
261 updated.

262

263 We also showed that we could lift-over nearly all protein-coding genes from GRCh38 to
264 the chimpanzee genome, PTRv2, with an average sequence identity of 98.7%. This is
265 consistent with previous work showing the human genome and chimpanzee genome
266 are approximately 98% identical (Chimpanzee Sequencing and Analysis Consortium,
267 2005). Comparing the gene order revealed 4 large regions on the homologs of
268 chromosomes 4, 5, 12, and 17 where the gene order is inverted. These regions are
269 consistent with previous reports: the chimpanzee genome has 9 well-characterized
270 pericentric inversions on chromosome homologs 1, 4, 5, 9, 12, 15, 16, 17 (Yunis and
271 Prakash, 1982). The 4 largest of these inversions are on 4, 5, 12, and 17 (Soto *et al.*,
272 2020) hence their visibility at this scale. Additionally, the co-linear mapping of genes
273 from human chromosome 2 to chimpanzee chromosomes 2A and 2B is consistent with
274 the known telomeric fusion of these chromosomes (Yunis and Prakash, 1982). The
275 consistency of the gene sequence identity with the known genome sequence identity
276 between chimpanzee and human, and the consistency of the gene order with the known
277 structural differences between the two genomes demonstrate the accuracy of Liftoff's
278 gene placements in a cross-species lift-over.

279
280 Annotating new assemblies with a lift-over strategy rather than *de novo* is limited in that
281 the annotation of the new assembly will only be as complete as the reference. However,
282 as reference annotations continue to improve through manual curation, experimental
283 validation or improved computational methods, Liftoff will enable easy integration of
284 these improvements across many genomes. We anticipate that Liftoff will be a valuable
285 tool in improving our understanding of the biological function of the large and rapidly
286 growing number of sequenced genomes.

287

288

289 **Acknowledgements**

290 This work was supported in part by NIH under grants R01-HG006677 and R35-
291 GM130151.

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306 **References**

- 307 Alonge, M. *et al.* (2020) Chromosome-scale assembly of the bread wheat genome,
308 *Triticum aestivum*, reveals over 5700 new genes. *bioRxiv*, 2020.04.06.028746.
- 309 Chimpanzee Sequencing and Analysis Consortium (2005) Initial sequence of the
310 chimpanzee genome and comparison with the human genome. *Nature*, **437**, 69–87.
- 311 Church, D.M. *et al.* (2011) Modernizing reference genome assemblies. *PLoS Biol.*, **9**,
312 e1001091.
- 313 Dijkstra, E.W. and Others (1959) A note on two problems in connexion with graphs.
314 *Numer. Math.*, **1**, 269–271.
- 315 Frankish, A. *et al.* (2019) GENCODE reference annotation for the human and mouse
316 genomes. *Nucleic Acids Res.*, **47**, D766–D773.
- 317 Gao, B. *et al.* (2018) segment_liftover : a Python tool to convert segments between
318 genome assemblies. *F1000Res.*, **7**, 319.
- 319 Guo, Y. *et al.* (2017) Improvements and impacts of GRCh38 human reference on high
320 throughput sequencing data analysis. *Genomics*, **109**, 83–90.
- 321 Harrow, J. *et al.* (2012) GENCODE: the reference human genome annotation for The
322 ENCODE Project. *Genome Res.*, **22**, 1760–1774.
- 323 He, Y. *et al.* (2019) Long-read assembly of the Chinese rhesus macaque genome and
324 identification of ape-specific structural variants. *Nat. Commun.*, **10**, 4233.
- 325 Jiao, Y. *et al.* (2017) Improved maize reference genome with single-molecule
326 technologies. *Nature*, **546**, 524–527.
- 327 Kuhn, R.M. *et al.* (2013) The UCSC genome browser and associated tools. *Brief.*
328 *Bioinform.*, **14**, 144–161.
- 329 Li, H. (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**,
330 3094–3100.
- 331 Shirley, M.D. *et al.* (2015) Efficient ‘pythonic’ access to FASTA files using pyfaidx PeerJ
332 PrePrints.
- 333 Shumate, A. *et al.* (2020) Assembly and annotation of an Ashkenazi human reference
334 genome. *Genome Biol.*, **21**, 129.
- 335 Soto, D.C. *et al.* (2020) Identification of Structural Variation in Chimpanzees Using
336 Optical Mapping and Nanopore Sequencing. *Genes*, **11**.
- 337 Yunis, J.J. and Prakash, O. (1982) The origin of man: a chromosomal pictorial legacy.
338 *Science*, **215**, 1525–1530.
- 339 Zhao, H. *et al.* (2014) CrossMap: a versatile tool for coordinate conversion between
340 genome assemblies. *Bioinformatics*, **30**, 1006–1007.