

1 Insertion-deletions are depleted in protein regions with predicted secondary structure

2
3 Yi Yang, Matthew Braga, and Matthew D. Dean

4
5 Molecular and Computational Biology

6 University of Southern California

7 1050 Childs Way

8 Los Angeles, CA 90089

9
10 matthew.dean@usc.edu

11

12

13

Abstract

14
15
16 A fundamental goal in evolutionary biology and population genetics is to understand
17 how selection shapes the fate of new mutations. Here we test the null hypothesis that
18 insertion-deletion events (indels) in protein coding regions occur randomly with respect
19 to secondary structures. We identified indels across 11,444 sequence alignments in
20 mouse, rat, human, chimp, and dog genomes, then quantified their overlap with four
21 different types of secondary structure – alpha helices, beta strands, protein bends, and
22 protein turns – predicted by deep-learning methods of AlphaFold2. Indels overlapped
23 secondary structures 54% as much as expected, and were especially under-
24 represented over beta strands, which tend to form internal, stable regions of proteins. In
25 contrast, indels were enriched by 155% over regions without any predicted secondary
26 structures. These skews were stronger in the rodent lineages compared to the primate
27 lineages, consistent with population genetic theory predicting that natural selection will
28 be more efficient in species with larger effective population sizes. Nonsynonymous
29 substitutions were also less common in regions of protein secondary structure, although
30 not as strongly reduced as in indels. In a complementary analysis of thousands of
31 human genomes, we showed that indels overlapping secondary structure segregated at
32 significantly lower frequency than indels outside of secondary structure. Taken together,
33 our study shows that indels are selected against if they overlap secondary structure,
34 presumably because they disrupt the tertiary structure and function of a protein.

35

36 Keywords: insertion, deletion, indels, evolution, selection

37

38

39

Significance

40 How do insertion-deletion mutations, which occur when short stretches of amino acids
41 are either added or deleted from a protein, accumulate in genomes? Here we show that
42 insertion-deletion events are less common in regions of proteins that are predicted to
43 form secondary structures. We present multiple lines of evidence to show that this is
44 most likely caused by selection against insertion-deletion events that disrupt secondary
45 structure, and therefore the overall function of a protein.

46

47

Introduction

48 Understanding the fate of new mutations is critical to defining the evolutionary

49 processes that shape biological diversity. At the level of single nucleotides, a rich body

50 of theory has been developed to infer whether mutations are neutral, deleterious, or

51 beneficial (Nielsen & Slatkin 2013; reviewed by Hedrick 2005; Hartl & Clark 2007).

52 Understanding the selective impact of insertion-deletion events (indels), which can

53 extend many nucleotides, has proven to be much more complicated.

54 Previous studies investigating the functional impact of indels generally fall into

55 two categories (Savino et al. 2022). First, protein engineering studies have shown that

56 indels can impact a protein's function, especially if they overlap important secondary

57 structures (Simm et al. 2007; Arpino et al. 2014; Tóth-Petróczy & Tawfik 2014; Gavrillov

58 et al. 2015; Grocholski et al. 2015; Liu et al. 2015, 2016; Jackson et al. 2017; Gavrillov

59 et al. 2018; Halliwell et al. 2018; Gonzalez et al. 2019; Woods et al. 2023). For example,

60 Liu et al. (2016) found that experimentally deleting amino acids in beta strands and

61 alpha helices of Green Fluorescent Protein tended to reduce fluorescence, while

62 deletions outside such regions were relatively neutral.

63 Second, evolutionary and population genetic studies have suggested that indels

64 are relatively deleterious if they are long (Pascarella & Argos 1992; Taylor et al. 2004;

65 Tao et al. 2007; Hsing & Cherkasov 2008; Kim & Guo 2010; Mills et al. 2011; Rockah-

66 Shmuel et al. 2013; Lek et al. 2016; Zhang et al. 2018), cause frame-shifts (Chen & Guo

67 2021; Montgomery et al. 2013; Chong et al. 2013; Iengar 2012; Bermejo-Das-Neves et

68 al. 2014), occur internally in the protein (Lin et al. 2017), alter flanking amino acids

69 (Zhang et al. 2011), or fall outside of disordered regions (Taylor et al. 2004; Bermejo-

70 Das-Neves et al. 2014; Khan et al. 2015; Light, Sagit, Ekman, et al. 2013; Light, Sagit,
71 Sachenkova, et al. 2013). Protein families with indels tend to diverge in their structure
72 and function relative to protein families without indels (Gavrilov et al. 2018, 2015;
73 Banerjee et al. 2019; Zhang et al. 2018, 2010; Jayaraman et al. 2022; Hormozdiari et al.
74 2009; Salari et al. 2008), suggesting indels can be an important source of evolutionary
75 novelty. Indeed, one study estimated that >70% of indels that have reached fixation
76 have done so through positive selection (Barton & Zeng 2019).

77 Two important evolutionary studies identified orthologs across species and then
78 overlapped inferred indels with experimentally determined protein structures in the
79 Protein Data Bank (PDB, Berman et al. 2000). Following the publication of the human,
80 mouse and rat genomes, Taylor et al. (2004) identified 52 orthologous protein-coding
81 genes that had an indel *and* a protein structure. Of these 52 indels, 31.5% of their
82 sequence overlapped secondary structure of any kind, compared to 52.5% expected. A
83 few years later, de la Chaux et al. (2007) analyzed the distribution of 343 protein-coding
84 indels identified from human-chimp-rhesus orthologs that also occurred in the PDB.
85 They found a deficiency of indels that overlapped alpha helices, but no difference in
86 indels that overlapped beta strands.

87 As impactful as these studies were, they may not paint a full picture of the
88 functional consequences of indel variation. The set of genes that could be studied was
89 small, mostly limited by structural protein data or annotated Pfam domains. Pfam
90 domains do not necessarily correlate with 3D structure and the PDB represents a
91 biased set of proteins (or protein regions) that are amenable to the experimental
92 approaches required for structural proteomics, such as their ability to be crystallized.

93 The relatively biased set of proteins for which we have structural data thus limits a
94 systematic analysis across full genomes. For example, one study of duplicated genes
95 could not analyze full-length proteins because of divergence between aligned gene
96 sequences and proteins represented in the PDB (Guo et al. 2012). However, the recent
97 release of AlphaFold2 – a deep-learning project that accurately predicts the 3D
98 structure of a protein from its amino acid sequence (Jumper et al. 2021; Varadi et al.
99 2022) – provides a unique opportunity to systematically study indels across full proteins
100 and whole genomes.

101 Here we combine genome-wide predictions of AlphaFold2 with evolutionary and
102 population genetic methods to ask whether indels occur randomly with respect to
103 secondary structure, providing the most comprehensive evolutionary investigation into
104 the fate of indels in protein coding regions. We report four main results: 1) 97,382 indels
105 identified from 11,444 five-species alignments in the tree (dog, ((mouse, rat), (human,
106 chimp))) overlapped secondary structures 54% as often as expected, but were 155%
107 more common than expected in regions with no predicted secondary structures, 2)
108 indels that overlapped beta strands and occurred internally in a protein were especially
109 rare, consistent with the known importance of these regions in overall protein structure,
110 3) skews in observed vs. expected were stronger in the rodent lineages compared to
111 the primate lineages, consistent with theory predicting more efficient selection in rodents
112 given their larger effective population sizes, and 4) within human populations, indels that
113 overlapped secondary structures occurred at significantly lower frequency compared to
114 indels outside of secondary structures. Taken together, our results indicate selection
115 acts against indels when they arise over structurally important regions of proteins,

116 presumably because they can disrupt overall structure and therefore the function of a
117 protein.

118

119

Materials and Methods

120 **Interspecific insertion-deletion (indel) events.** We downloaded protein sequences
121 from all protein-coding genes identified as one-to-one orthologs between mouse, rat,
122 human, chimp, and dog from Ensembl version 107 (ensembl.org). In the case of
123 alternative transcripts, we chose the longest translated transcript to represent the gene.
124 11,444 genes had one-to-one orthologs across all five species.

125 We aligned proteins using GUIDANCE (Penn, Privman, Landan, et al. 2010; Penn,
126 Privman, Ashkenazy, et al. 2010; Privman et al. 2012; Levy Karin et al. 2014). This
127 approach estimates per-site alignment confidence by calculating its consistency across
128 different starting guide trees, allowing us to incorporate a measure of confidence in
129 downstream analyses. Importantly, we could use GUIDANCE scores to estimate error in
130 indel placement and identify indels that were confidently placed. In each GUIDANCE
131 iteration, we aligned protein sequences with MAFFT (Kato et al. 2002). We ran MAFFT
132 under the recommended default parameters; in the case of indels the most important
133 default parameters were the gap opening penalty (default=1.53) and gap offset value
134 (similar to gap extension penalty, default=0.123). We then identified all indels as gaps
135 from all 11,444 alignments (Fig. 1).

136 Our analyses could be impacted by sequencing errors or annotation errors that
137 result in spurious inclusion or exclusion of amino acids from certain genes, or by
138 alignment errors (Chowdhury & Garai 2017; Fitch & Smith 1983). Therefore, we

139 repeated all downstream analyses after subsetting indels in four different ways: 1)
140 INTERNAL: any indels that reached the beginning or ends of alignments were excluded,
141 as visual inspection indicated these were noisy regions of alignment that could be
142 related to incomplete annotation of full length genes, 2) GU94 PA100 GD40:
143 INTERNAL indels whose flanking five positions on both 5' and 3' ends (10 flanking
144 positions total) had an average GUIDANCE confidence score of at least 0.94 (median
145 observed), contained no overlapping indels, and had an average Grantham distance
146 (Grantham 1974) of less than 40 (median observed), where Grantham distance was
147 calculated using the R package AGVGD (<https://CRAN.R-project.org/package=agvgd>).
148 This subset was meant to enrich for well-anchored indels and avoid problems
149 distinguishing gaps in alignment due to protein divergence, versus gaps in alignment to
150 insertion-deletion events (Salari et al. 2008; Jilani et al. 2022; Snir & Pachter 2006), 3)
151 LENGTH LTE20: INTERNAL indels that were less than or equal to 20 amino acids long
152 in length, minimizing the impact of large indels that sometimes appeared to be spurious,
153 and 4) MERGED: INTERNAL indels after merging coordinates that overlapped, so that
154 sites in an alignment that were in different overlapping regions only contributed once.
155 We present the results from these four subsets as supplementary files, but they all
156 produced essentially identical results as analyzing ALL indels.

157

158 **AlphaFold2.** AlphaFold2 is a deep learning approach developed by DeepMind to
159 predict the 3D structure of proteins from only their amino acid sequence (Jumper et al.
160 2021; Varadi et al. 2022). Comparison to empirical data indicates these computational
161 predictions are over 90% accurate.

162 AlphaFold2 assigns 43 different secondary structures to different regions of a
163 protein, which we collapsed into five main categories. There were 32 different
164 AlphaFold2 predictions that contained the phrase HELX, which are predictions of
165 different helices; we collapsed these into the single term HELIX. There were 8 different
166 AlphaFold2 predictions that contain the phrase TURN, which are regions where the
167 polypeptide is predicted to reverse direction in 3D space; we collapsed these into the
168 single term TURN. We included the single AlphaFold2 prediction STRAND as-is, which
169 are regions predicted to contain beta strands (also referred to as beta sheets). We
170 included the single AlphaFold2 prediction BEND as-is, which are regions where the
171 polypeptide is predicted to change direction but not fully reverse. There was one last
172 AlphaFold2 prediction OTHER, but we did not observe any instances of this prediction in
173 any of the proteins analyzed in this study so ignored that term. Each residue in the
174 Uniprot protein used by AlphaFold2 was assigned to one of these four categories, or
175 assigned the term NONE if they occurred outside any predicted secondary structure.

176 To link AlphaFold2 predictions to our five-species alignments above, we included
177 the Uniprot sequence in the alignment (Fig. 1). In rare cases, the AlphaFold2-
178 downloaded Uniprot sequence did not match the Ensembl-downloaded Uniprot
179 sequence, in which case we discarded the alignment from all analyses. Each position in
180 each indel was then assigned HELIX, STRAND, TURN, BEND, or NONE (Fig. 1). In
181 cases where the Uniprot sequence was “deleted” (for example, indel 50-52 in Fig. 1),
182 we assigned one-half of the deleted positions to whatever was assigned to its 5'-flanking
183 residue, and one-half to whatever was assigned to its 3'-flanking residue.

184

185 **Randomization of indel positions.** We generated null expectations through a
186 randomization procedure. For each alignment, we randomly shuffled the starting
187 position of each indel, then extended each randomized indel by its observed length. In
188 cases where a randomized indel extended past the end of an alignment, we wrapped
189 the randomized indel to the front of the alignment. After shuffling the unique indels
190 within each alignment, we re-calculated the number of residues falling in each
191 secondary structure, exactly as described above. We repeated this process 200 times
192 to generate null expectations. We repeated this entire process for the four different
193 subsets described above. For these four subsets, the relevant alignments were first
194 truncated to match included regions and provide a more appropriate background for
195 randomization.

196

197 **Gene Ontology enrichment.** For the MERGED indels only, we identified relative
198 outliers by counting the number of sites in the alignment overlapping NONE vs. not,
199 versus sites overlapping indels vs. not. We excluded alignments that had fewer than 5
200 positions in any of these four cells of this 2x2 table, then applied a X^2 test and corrected
201 resulting p-values (Benjamini & Hochberg 1995). Genes with a $-\log_{10}$ p.value of at least
202 10 and at least a 1.5 fold change in expectation were taken as relative outliers. We
203 tested whether these relative outlier genes were enriched for any Biological Process,
204 Molecular Function, or Cellular Component using Panther Classification system
205 (Thomas et al. 2022; Mi et al. 2017, 2019, 2013), run from PantherDB
206 (<https://pantherdb.org/>), with the settings “Test Type=Fisher’s Exact Test” and

207 “Correction=Calculate False Discovery Rate”. We also performed Gene Ontology
208 analyses for genes which had no indels across the five species analyzed.

209

210 **Accessibility and pIIDD scores.** Sites that are relatively internal on a 3D
211 protein evolve more slowly than external sites, both at the level of nonsynonymous
212 mutations (Tóth-Petróczy & Tawfik 2011; Goldman et al. 1998; Shahmoradi et al. 2014;
213 Franzosa & Xia 2009; Shih et al. 2012; Yeh et al. 2014; Dean et al. 2002; Scherrer et al.
214 2012; Bustamante et al. 2000; Marsh & Teichmann 2014) and indel variation (Guo et al.
215 2012; Hsing & Cherkasov 2008). This correlation is complicated by whether or not
216 external residues interact with other proteins (Kim et al. 2006; Mintseris & Weng 2005),
217 or if externally oriented residues form active sites of proteins (Slodkowicz & Goldman
218 2020). For each site in each alignment, we calculated relative solvent accessibility,
219 which is the degree to which a residue occurs on the outside of a folded protein (Tien et
220 al. 2013), using FREESASA (Mitternacht 2016) with the “--format=rsa” option, using the
221 AlphaFold2 structure as input. We also compared pIIDD scores (Mariani et al. 2013)
222 across an alignment. pIIDD scores are computational measures of confidence included
223 in AlphaFold2 predictions. According to AlphaFold2, pIIDD scores <50 likely represent
224 intrinsically disordered or unstructured regions. As above, any “deletions” in the Uniprot
225 sequence were divided, and one-half of their sites were assigned the accessibility and
226 pIIDD scores of their 5' flanking residue, and the other half to the scores of their 3'
227 flanking residue.

228 As will be shown below, secondary structure and relative solvent accessibility are
229 strongly correlated. In an attempt to separate the effects of these two features on the

230 probability of observing an indel, we compared Receiver Operating Characteristic
231 (ROC) curves and Area Under the Curve (AUC) values from three Generalized Linear
232 Models and then compared their likelihoods. Two models tested whether the probability
233 of observing an indel was a function of secondary structure or relative solvent
234 accessibility alone – $\text{glm}(\text{indel} \sim \text{secondary_structure})$ or $\text{glm}(\text{indel} \sim \text{rsa})$, respectively. A
235 third model included both as independent variables – $\text{glm}(\text{indel} \sim \text{secondary_structure} +$
236 $\text{rsa})$. We quantified the gain in likelihood when we included both independent variables,
237 versus each one separately. For all three models we included the “family = binomial”
238 argument to model logistic variance. Our approach closely followed that of Jackson et
239 al. (2017), modifying their scripts to suit our approach.

240 Because sites in a protein are not independent from each other, before applying
241 Generalized Linear Models we randomly sampled a single site from each alignment.
242 However, we did not sample sites with equal probability. Instead, we downweighted the
243 probability of sampling by the inverse of the grand total of the five secondary structures
244 (HELIX, STRAND, TURN, BEND, or NONE). By including this weighting scheme, we
245 ensured even sampling of secondary structures, increasing power of all three
246 Generalized Linear Models.

247

248 **Comparison to synonymous and nonsynonymous mutations.** To provide additional
249 context with which to interpret the distribution of indels, we tested three different
250 nucleotide-based sites. First, we quantified the distribution of invariant sites across
251 secondary structure as a kind of null distribution. Then we quantified the same with
252 respect to synonymous and nonsynonymous sites. We predicted that synonymous sites

253 should distribute similarly to invariant sites, because they do not alter the protein
254 sequence and thus probably have relatively minor effect on secondary structure.
255 Conversely, we predicted that nonsynonymous sites would occur less frequently over
256 secondary structure because, all else equal, their resulting amino acid changes could
257 alter secondary structure.

258 Using the same 5-species alignments above, we reverse-translated each protein
259 to its transcript, downloaded from Ensembl version 107. We counted the proportion of
260 synonymous vs. nonsynonymous variants occurring over the different secondary
261 structures, compared to invariant sites. We only quantified synonymous vs.
262 nonsynonymous variants from the same alignments and sites that were used in our
263 indel analyses.

264

265 **Intraspecific indel events.** As a complementary analysis to the interspecific analyses
266 described above, we analyzed intraspecific variation from Phase 3 of the 1000 Human
267 Genomes project ([https://www.internationalgenome.org/data-portal/data-collection/30x-](https://www.internationalgenome.org/data-portal/data-collection/30x-grch38)
268 [grch38](https://www.internationalgenome.org/data-portal/data-collection/30x-grch38)) (The Genomes Project 2015; Byrska-Bishop et al. 2022). This database
269 contains haplotype-phased indel calls (files named like
270 ALL.chr1.shapeit2_integrated_snvindels_v2a_27022019.GRCh38.phased.INDELS.vcf)
271 from 2,504 unrelated samples from 26 populations, with sample size ranging from 61 to
272 113 per population. These 26 populations derive from five large geographic areas:
273 Africa, East Asia, South Asia, South America, and Europe.

274 Indel coordinates were truncated to match exon coordinates downloaded from
275 UCSC Table Browser (table name=unipAliSwissprot from GRCH38). For any protein-

276 coding genes that contained at least one indel, we assembled the reference and
277 alternative alleles from the human genome, computationally placed indels, and then
278 translated both alleles. Any indels that resulted in a frameshift in the first 95% of the
279 protein-coding transcript (counted from 5' translation start site) were excluded, because
280 it is unclear whether reference and alternative alleles share 3D structure if they are
281 dramatically frame-shifted with respect to each other.

282 We only analyzed genes that were part of the five-species interspecific analyses
283 described above. Otherwise, we would have included recent human-specific duplicates,
284 where predictions might become noisy because of uncertainty about the exact timing of
285 duplication along the lineage to modern humans.

286

287

Results

288 **Indels were depleted in regions with secondary structure.** There were 11,444
289 genes that had one-to-one orthologs between dog, mouse, rat, chimp, and human
290 genomes. Across these 11,444 alignments we identified 97,382 indels spanning
291 1,272,048 positions. Indel sizes ranged from 1 to 2,870 residues long, but most were
292 small: the 25%, 50%, and 75% quantiles were 1, 3, and 10 residues, respectively. Indel
293 positions overlapped secondary structures significantly less than expected (Fig. 2, Table
294 1). Indel positions were most under-represented in STRAND, occurring at 43%
295 expectations (calculated as 55,293 indel sites that overlapped STRAND, compared to
296 129,070 averaged across 200 randomizations), followed by indel positions occurring in
297 TURN (55%), HELIX (57%), and BEND (59%) (Table 1). In contrast, indel positions
298 occurred at 155% expectation in NONE, meaning indels were much more likely occur in

299 protein regions with no predicted secondary structure (Table 1). All observed values fell
300 far outside the distributions from randomization (Fig. 2), translating into a p-value of
301 essentially 0. We reached nearly identical conclusions after subsetting indels in four
302 ways described above (Supplementary Figure 1, Supplementary Table 1), with one
303 exception: indels over TURN and BEND are not under-represented in the very stringent
304 subset GU94 PA100 GD40 (Supplementary Figure 1, Supplementary Table 1).

305

306 **Skews in indel distribution were stronger in rodents.** By using dog as an outgroup,
307 we polarized all indels into either an insertion or deletion and placed each indel event on
308 a specific branch in the phylogenetic tree, using simple parsimony. In other words, if
309 amino acid sequences existed for mouse and rat, but not for the other species, that
310 indel was mapped as an insertion on the branch leading to rodents.

311 There are seven branches on the phylogenetic tree analyzed here. Across the
312 four secondary structures (BEND, TURN, STRAND, and HELIX), 24 of 28 O:E values
313 were lower for insertions compared to deletions (Figure 3). Conversely, across NONE
314 sites all branches showed higher O:E for insertions compared to deletions. Taken
315 together, these results suggest that insertions over secondary structure are more
316 deleterious than deletions.

317 For the four secondary structures, O:E values were consistently lower in rodent
318 lineages compared to primate lineages. There are four secondary structure that can be
319 mapped to three rodent branches and three primate branches, where each branch
320 contains insertions and deletions, for a total of 48 O:E values in Figure 2. 46 of these 48
321 O:E values were lower in the rodent lineages compared to primate lineages. For

322 example, O:E values for insertions over STRAND in the three rodent lineages = 0.26,
323 0.39, and 0.35, while in primates the three values = 0.52, 0.46, and 0.41. Conversely,
324 O:E values for NONE sites tend to be higher in rodents compared to primates. In sum,
325 indels were especially unlikely to overlap secondary structures in rodents. All patterns
326 described held after analyzing the four different subsets of indels described above
327 (Supplementary Figure 2).

328

329 **GO analysis.** We identified 797 alignments (genes) where the enrichment of indels over
330 NONE was especially high. Compared to the rest of the 4,995 alignments, these 797
331 genes showed no statistical enrichment of Biological Process, but under the Cellular
332 Component and Molecular Function ontologies showed enrichment of terms associated
333 with cilia and ubiquitination. This enrichment lacks an obvious explanation.

334 We identified 88 alignments (genes) whose indels overlapped NONE much less
335 than expected. None of these 88 genes showed enrichment of Biological Process or
336 Molecular Function but showed enrichment of gene products localized to the nucleus
337 under Cellular Component. In sum, there were no striking or consistent patterns of
338 Gene Ontology enrichment associated with outlier genes in either direction.

339 We also analyzed the 904 genes which had no indels across any of the five
340 species in the alignment. GO analysis uncovered many functional terms associated with
341 neurotransmission, including synapse localization and synaptic transmission
342 (Supplementary Table 2). This result suggests that genes involved in neurotransmission
343 may be especially intolerant of indel mutations. Interestingly, genes involved in immune

344 response appeared to be under-represented among genes with no indels. This result
345 may indicate that immune genes undergo indel mutations more often than expected.

346

347 **Indels were enriched in regions with high accessibility and low pI/DDT scores.**

348 Accessibility and pI/DDT scores varied according to secondary structure. STRAND had
349 low accessibility and high pI/DDT scores, indicating these secondary structures tend to
350 fall on the inside of proteins and are relatively stable (Fig. 4). On the other end of the
351 spectrum, NONE sites were much more accessible, with lower pI/DDT scores, indicating
352 external and unstable regions of proteins (Fig. 4).

353 Importantly, sites that overlapped indels consistently showed higher accessibility
354 and lower pI/DDT scores (compare X vs. O within each group, Fig. 4). In other words,
355 *within* each secondary structure, indels were more commonly observed at sites that
356 were relatively external and in relatively unstable regions, compared to sites that did not
357 overlap indels. Woods et al. (2023) found that experimentally deleting amino acids that
358 reside in regions of high pI/DDT were most likely to have a deleterious effect on protein
359 function, providing an explanation for why we observe indels more frequently in regions
360 with low pI/DDT scores. This pattern held across all four subsets of indels described
361 above (Supplementary Figure 3).

362 Comparing three different Generalized Linear Models demonstrated that the
363 effects of secondary structure were indistinguishable from the effects of relative solvent
364 accessibility (Table 2). In the ALL dataset, secondary structure performed about as well
365 as relative solvent accessibility (AUC=0.684 vs. 0.707, respectively), and including both
366 as independent variables had only minor improvement to AUC (0.720) compared to

367 single regressions. Similar results were obtained across the four subsets of data
368 described above (Table 2). This shows that secondary structure and relative solvent
369 accessibility are so correlated with each other that their effects cannot be meaningfully
370 separated.

371

372 **Nonsynonymous variants were also depleted in protein regions with secondary**
373 **structure.** Among the 11,444 alignments, we analyzed 3.8, 2.14, and 1.67 million
374 codons that were invariant, synonymous, or nonsynonymous, respectively (Table 1).
375 Synonymous codons overlapped secondary structures as often as invariant codons
376 (synonymous-to-invariant ratios ranging from 0.86 to 1.17, Table 1). In contrast,
377 nonsynonymous codons occurred far less frequently across the four secondary
378 structures (nonsynonymous-to-invariant ratios ranging from 0.71 to 0.92) and more over
379 NONE (nonsynonymous-to-invariant ratio of 1.24) (Table 1). These nonsynonymous-to-
380 invariant ratios were generally smaller in magnitude than the O:E ratios estimated from
381 indel distribution (Table 1). For example, indels occurred at 43% expectation over
382 STRAND, while nonsynonymous codons occurred at 71% “expectation” (Table 1).

383 Similar patterns emerged after analyzing the four subsets of indels
384 (Supplementary Table 1). The main exception was that nonsynonymous-to-invariant
385 ratios ranged from 0.91 to 0.98 across the four secondary structures, and from 1.05 to
386 1.09 for NONE (Supplementary Table 1). In other words, we still observed the general
387 pattern that nonsynonymous variants were under-represented across the four
388 secondary structures and enriched over NONE, although at a smaller magnitude
389 compared to the overall analysis.

390

391 **Human intraspecific variation.** We identified 1,921 indels from 1,436 unique genes,
392 comprising a total of 4,354 positions. Most of these occurred at a frequency of 1 allele
393 observed among 5,008 phased alleles in the 1000 genomes project. We did not exclude
394 these; even if they are due to sequencing or mapping errors, there is no reason to
395 believe they would inflate our overall false positive rate as such errors should occur
396 blindly with respect to secondary structure of proteins. In addition, an indel at a
397 frequency of 1 allele could be especially deleterious, so we included them.

398 Across all 6 geographic regions, indel sites spanning NONE occurred at nearly
399 twice the frequency than secondary structures. NONE indels reached a mean frequency
400 of 4 alleles out of 5,008 phased alleles, compared to BEND/HELIX/TURN indels (3
401 alleles) and STRAND indels (1 allele) (Kruskal-Wallis $X^2= 37.8$, $df = 2$, $p\text{-value} < 10^{-8}$). If
402 we use a minor allele frequency cutoff of 1%, 3% or 5% these patterns disappear,
403 indicating that the majority of signal comes from the fact that a large proportion of
404 STRAND indels occur as singletons.

405

406

Discussion

407 Our study combined the recent revolution in protein structure, ushered in by the
408 AlphaFold2 project (Jumper et al. 2021), with evolutionary, population genetic, and
409 permutation-based analyses to demonstrate that indels were depleted in regions of
410 predicted secondary structure. This skew is especially strong for STRAND, which is
411 consistent with these structures being internal and stable regions that are important for
412 the overall 3D structure of a protein (Echave et al. 2016).

413 There are two non-mutually exclusive models – a mutational bias model versus a
414 selection model – that could explain the non-random distribution of indels that we
415 observe here. Under a mutational bias model, the four secondary structures experience
416 fundamentally different rates of indel mutation. The four different secondary structures
417 tested here display systematic differences in amino acid composition (Chou & Fasman
418 1975; Fujiwara et al. 2012), which predicts different base composition and/or repetitive
419 elements in the underlying DNA, which in turn could influence mutation rate.

420 However, three patterns in our data argue against the mutational bias
421 hypothesis, and instead provide support for a model where selection acts against indels
422 that are more likely to disrupt protein function. First, *within* each secondary structure,
423 positions with indels tend to occur in externally oriented and high-pIDDT regions of
424 proteins (Fig. 4). A mutational bias hypothesis cannot account for this discrepancy
425 because they are the same secondary structures in different parts of the same protein.
426 Second, the observed vs. expected ratios (Table 1) are stronger in rodents compared to
427 primates (Figure 3). A mutational bias hypothesis cannot account for this interspecific
428 variation unless different species also experience different mutational biases. In
429 contrast, this pattern is predicted by a model of selection, because natural selection will
430 operate more efficiently in species with large effective population size (Lynch 2007;
431 Kimura 1983; Charlesworth 2009). Rodents have an effective population size that is
432 roughly 10-fold larger than primates (Geraldes et al. 2011; Zhao et al. 2000; Won & Hey
433 2005; Ohta 1972; Geraldes et al. 2008). Finally, we showed that nonsynonymous
434 variants were also depleted in regions of secondary structure, although not to the same
435 degree (Table 1). A mutational bias hypothesis cannot explain the depletion of both

436 indels and nonsynonymous variants over secondary structure, because these two
437 classes differ in their mutational process.

438 To be sure, it is unlikely that indel mutations arise randomly. For example, G+C
439 content often correlates with a genomic region's susceptibility to insertions or deletions
440 (Sinden et al. 2002; Taylor et al. 2004), as well as features suggestive of a slippage
441 mechanism (Nishizawa & Nishizawa 2002). However, a model of selection does not
442 require indel mutation to be completely random. A selection model only requires any
443 non-randomness in mutational process to be equally distributed across the five
444 categories of secondary structure tested here. It should also be pointed out that our
445 study reports average deviations in observed vs. expected across the entire genome. It
446 remains unknown how much the strength of selection varies across individual indels,
447 although our Gene Ontology results did not uncover any functional similarity among the
448 most highly skewed genes.

449 It is noteworthy that even within humans, we observed proportionately fewest
450 indels over STRAND – exactly the secondary structure where indels were depleted in
451 our five species analyses. The low historical effective population size of humans,
452 coupled with multiple bottlenecks, are expected to reduce the efficiency of selection, yet
453 we still observe skews in indel locations.

454 In conclusion, our analyses indicate that any change in amino acid sequence is
455 likely to be deleterious for secondary structure, especially if that change is not a single
456 nonsynonymous mutation, but the insertion or deletion of multiple amino acids. Indels
457 that overlap STRAND and/or buried regions of the protein, appear to be the most
458 deleterious, while indels over NONE the least. By analyzing the AlphaFold2 predictions,

459 we have quantified these effects over whole genomes and full-length proteins, revealing
460 a role for protein structure on the evolution of its primary sequence.

461

462 **Data and resource availability**

463 All data, code, and intermediate files required to reproduce the results here, as well as a
464 README file, are available on Dryad (<https://doi.org/10.5061/dryad.bk3j9kdk9>) as a
465 single protein_structure.tar.gz file (8.5 Gb). [for reviewers only: that link is not yet public;
466 this link provides access:

467 <https://datadryad.org/stash/share/5NLwY6lUt75oIqY16DqFRkywjBUx2eoela6RYDHFHdg>
468 g]

469 **Acknowledgements**

470 Cornelius Gati provided many useful discussions on protein structure. Mark Chaisson
471 provided many useful suggestions regarding 1000 genomes data. Jeff Jensen
472 discussed population genetic inference. We thank Jackson et al. (2017) for publishing
473 their supplementary data; our Generalized Linear Modeling largely borrowed their code.
474 Two anonymous reviewers and an associate editor greatly improved the quality of the
475 manuscript. This project was born in BISC444; we thank the students of that class for
476 their input. Computation was supported by Center for Advanced Research Computing
477 (CARC) at the University of Southern California – Tomek Osinski provided many useful
478 discussions. Funding was provided by National Science Foundation grant #2027373.

479 **Figure Legends**

480

481 **Figure 1.** Schematic of main methodology. Shown is a hypothetical protein alignment
482 between five species, which identified two unique indel events (positions 50-52 and
483 positions 530-534). By including the Uniprot sequence from AlphaFold2, we mapped
484 from indel coordinates into predicted secondary structures. In this example, three
485 positions fell over HELIX and five positions fell over SHEET. During randomization, we
486 would permute the starting locations of these two indel events, then extend them by
487 their observed length. Intraspecific analyses of human genomes proceeded in almost
488 the same manner, except that indels were already called in their corresponding .vcf
489 files.

490

491 **Figure 2.** Comparison of observed vs. expected number of alignment positions that
492 overlap indels in the 11,444 alignments, stratified by secondary structure. Histograms
493 built from randomizing indel positions across the alignments. Arrows at top originate at
494 the mean expectation for each group, and terminate at the observed value. Indel sites
495 overlap NONE 132% more than expected, and overlap the four secondary structures
496 less than expected (ranging from 62% expectation in STRAND to 84% expectation in
497 TURN). Also see Table 1.

498

499 **Figure 3.** Observed:Expected ratios of indels, polarized into insertions (above branch)
500 versus deletions (below branch), using Dog as outgroup. There is no consistent

501 difference in O:E in insertions and deletions, but the branches leading to rodent species
502 generally show stronger skews than branches leading to primates.

503

504 **Figure 4.** Weighted means of relative solvent accessibility (red, left axis) and pI-DDT
505 scores (blue, right axis) across secondary structures, stratified by sites occurring over
506 indels (X) versus sites not overlapping indels (O). Numbers on x axis indicate the
507 number of sites that overlap an indel versus not (separated by |).

508

509 **Figure 5.** Violin plot of the minor allele frequency of indels in protein coding regions,
510 segregating within humans, stratified by secondary structure. B/H/T = pooled

511 BEND+HELIX+TURN. Numbers on x-axis indicate number of positions observed.

512 Figure includes all human populations pooled; results remain qualitatively the same if
513 we analyze populations separately.

514

515

516

References

- 517
518
- 519 Arpino JAJ, Reddington SC, Halliwell LM, Rizkallah PJ, Jones DD. 2014. Random
520 Single Amino Acid Deletion Sampling Unveils Structural Tolerance and the
521 Benefits of Helical Registry Shift on GFP Folding and Structure. *Structure*.
522 22:889–898. doi: 10.1016/j.str.2014.03.014.
- 523 Banerjee A, Levy Y, Mitra P. 2019. Analyzing Change in Protein Stability Associated
524 with Single Point Deletions in a Newly Defined Protein Structure Database. *J.*
525 *Proteome Res.* 18:1402–1410. doi: 10.1021/acs.jproteome.9b00048.
- 526 Barton HJ, Zeng K. 2019. The Impact of Natural Selection on Short Insertion and
527 Deletion Variation in the Great Tit Genome. *Genome Biology and Evolution*.
528 11:1514–1524. doi: 10.1093/gbe/evz068.
- 529 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and
530 powerful approach to multiple testing. *Journal of the royal statistical society.*
531 *Series B (Methodological)*. 57:289–300.
- 532 Berman HM et al. 2000. The Protein Data Bank. *Nucleic Acids Research*. 28:235–242.
533 doi: 10.1093/nar/28.1.235.
- 534 Bermejo-Das-Neves C, Nguyen H-N, Poch O, Thompson JD. 2014. A comprehensive
535 study of small non-frameshift insertions/deletions in proteins and prediction of
536 their phenotypic effects by a machine learning method (KD4i). *BMC*
537 *Bioinformatics*. 15:111. doi: 10.1186/1471-2105-15-111.
- 538 Bustamante CD, Townsend JP, Hartl DL. 2000. Solvent Accessibility and Purifying
539 Selection Within Proteins of *Escherichia coli* and *Salmonella enterica*. *Molecular*
540 *Biology and Evolution*. 17:301–308. doi:
541 10.1093/oxfordjournals.molbev.a026310.
- 542 Byrska-Bishop M et al. 2022. High-coverage whole-genome sequencing of the
543 expanded 1000 Genomes Project cohort including 602 trios. *Cell*. 185:3426-
544 3440.e19. doi: 10.1016/j.cell.2022.08.004.
- 545 Charlesworth B. 2009. Effective population size and patterns of molecular evolution and
546 variation. *Nat Rev Genet*. 10:195–205.
- 547 de la Chaux N, Messer PW, Arndt PF. 2007. DNA indels in coding regions reveal
548 selective constraints on protein evolution in the human lineage. *BMC Evol Biol*.
549 7:191. doi: 10.1186/1471-2148-7-191.
- 550 Chen J, Guo J. 2021. Structural and functional analysis of somatic coding and UTR
551 indels in breast and lung cancer genomes. *Sci Rep*. 11:21178. doi:
552 10.1038/s41598-021-00583-1.

553 Chong Z et al. 2013. The Evolution of Small Insertions and Deletions in the Coding
554 Genes of *Drosophila melanogaster*. *Molecular Biology and Evolution*. 30:2699–
555 2708. doi: 10.1093/molbev/mst167.

556 Chou PY, Fasman GD. 1975. Conformational parameters for amino acids in helical, β -
557 sheet, and random coil regions calculated from proteins. ACS Publications. doi:
558 10.1021/bi00699a001.

559 Chowdhury B, Garai G. 2017. A review on multiple sequence alignment from the
560 perspective of genetic algorithm. *Genomics*. 109:419–431. doi:
561 10.1016/j.ygeno.2017.06.007.

562 Dean AM, Neuhauser C, Grenier E, Golding GB. 2002. The Pattern of Amino Acid
563 Replacements in α/β -Barrels. *Molecular Biology and Evolution*. 19:1846–1864.
564 doi: 10.1093/oxfordjournals.molbev.a004009.

565 Echave J, Spielman SJ, Wilke CO. 2016. Causes of evolutionary rate variation among
566 protein sites. *Nat Rev Genet*. 17:109–121. doi: 10.1038/nrg.2015.18.

567 Fitch WM, Smith TF. 1983. Optimal sequence alignments. *Proceedings of the National*
568 *Academy of Sciences*. 80:1382–1386. doi: 10.1073/pnas.80.5.1382.

569 Franzosa EA, Xia Y. 2009. Structural Determinants of Protein Evolution Are Context-
570 Sensitive at the Residue Level. *Molecular Biology and Evolution*. 26:2387–2395.
571 doi: 10.1093/molbev/msp146.

572 Fujiwara K, Toda H, Ikeguchi M. 2012. Dependence of α -helical and β -sheet amino acid
573 propensities on the overall protein fold type. *BMC Structural Biology*. 12:18. doi:
574 10.1186/1472-6807-12-18.

575 Gavrilov Y, Dagan S, Levy Y. 2015. Shortening a loop can increase protein native state
576 entropy. *Proteins: Structure, Function, and Bioinformatics*. 83:2137–2146. doi:
577 10.1002/prot.24926.

578 Gavrilov Y, Dagan S, Reich Z, Scherf T, Levy Y. 2018. An NMR Confirmation for
579 Increased Folded State Entropy Following Loop Truncation. *J. Phys. Chem. B*.
580 122:10855–10860. doi: 10.1021/acs.jpcc.8b09658.

581 Geraldes A et al. 2008. Inferring the history of speciation in house mice from autosomal,
582 X-linked, Y-linked and mitochondrial genes. *Mol Ecol*. 17:5349–63.

583 Geraldes A, Basset P, Smith KL, Nachman MW. 2011. Higher differentiation among
584 subspecies of the house mouse (*Mus musculus*) in genomic regions with low
585 recombination. *Mol Ecol*. 20:4722–36. doi: 10.1111/j.1365-294X.2011.05285.x.

586 Goldman N, Thorne JL, Jones DT. 1998. Assessing the Impact of Secondary Structure
587 and Solvent Accessibility on Protein Evolution. *Genetics*. 149:445–458. doi:
588 10.1093/genetics/149.1.445.

- 589 Gonzalez CE, Roberts P, Ostermeier M. 2019. Fitness Effects of Single Amino Acid
590 Insertions and Deletions in TEM-1 β -Lactamase. *Journal of Molecular Biology*.
591 431:2320–2330. doi: 10.1016/j.jmb.2019.04.030.
- 592 Grantham R. 1974. Amino acid difference formula to help explain protein evolution.
593 *Science*. 185:862–864.
- 594 Grocholski T, Dinis P, Niiranen L, Niemi J, Metsä-Ketelä M. 2015. Divergent evolution of
595 an atypical S-adenosyl-L-methionine–dependent monooxygenase involved in
596 anthracycline biosynthesis. *Proceedings of the National Academy of Sciences*.
597 112:9866–9871. doi: 10.1073/pnas.1501765112.
- 598 Guo B, Zou M, Wagner A. 2012. Pervasive Indels and Their Evolutionary Dynamics
599 after the Fish-Specific Genome Duplication. *Molecular Biology and Evolution*.
600 29:3005–3022. doi: 10.1093/molbev/mss108.
- 601 Halliwell LM et al. 2018. Δ Flucs: Brighter *Photinus pyralis* firefly luciferases identified by
602 surveying consecutive single amino acid deletion mutations in a thermostable
603 variant. *Biotechnology and Bioengineering*. 115:50–59. doi: 10.1002/bit.26451.
- 604 Hartl DL, Clark AG. 2007. *Principles of population genetics*. 4th ed. Sinauer:
605 Sunderland, MA.
- 606 Hedrick PW. 2005. *Genetics of populations*. 3rd ed. Jones and Bartlett: Boston.
- 607 Hormozdiari F et al. 2009. The Effect of Insertions and Deletions on Wirings in Protein-
608 Protein Interaction Networks: A Large-Scale Study. *Journal of Computational
609 Biology*. 16:159–167. doi: 10.1089/cmb.2008.03TT.
- 610 Hsing M, Cherkasov A. 2008. Indel PDB: A database of structural insertions and
611 deletions derived from sequence alignments of closely related proteins. *BMC
612 Bioinformatics*. 9:293. doi: 10.1186/1471-2105-9-293.
- 613 Iengar P. 2012. An analysis of substitution, deletion and insertion mutations in cancer
614 genes. *Nucleic Acids Research*. 40:6401–6413. doi: 10.1093/nar/gks290.
- 615 Jackson EL, Spielman SJ, Wilke CO. 2017. Computational prediction of the tolerance to
616 amino-acid deletion in green-fluorescent protein. *PLOS ONE*. 12:e0164905. doi:
617 10.1371/journal.pone.0164905.
- 618 Jayaraman V, Toledo-Patiño S, Noda-García L, Laurino P. 2022. Mechanisms of protein
619 evolution. *Protein Science*. 31:e4362. doi: 10.1002/pro.4362.
- 620 Jilani M, Haspel N, Jagodzinski F. 2022. Detection and Analysis of Amino Acid
621 Insertions and Deletions. In: *Algorithms and Methods in Structural Bioinformatics*.
622 Haspel, N, Jagodzinski, F, & Molloy, K, editors. Computational Biology Springer
623 International Publishing: Cham pp. 89–99. doi: 10.1007/978-3-031-05914-8_5.

- 624 Jumper J et al. 2021. Highly accurate protein structure prediction with AlphaFold.
625 Nature. 596:583–589.
- 626 Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple
627 sequence alignment based on fast Fourier transform. Nucleic acids research.
628 30:3059–3066.
- 629 Khan T, Douglas GM, Patel P, Nguyen Ba AN, Moses AM. 2015. Polymorphism
630 Analysis Reveals Reduced Negative Selection and Elevated Rate of Insertions
631 and Deletions in Intrinsically Disordered Protein Regions. Genome Biology and
632 Evolution. 7:1815–1826. doi: 10.1093/gbe/evv105.
- 633 Kim PM, Lu LJ, Xia Y, Gerstein MB. 2006. Relating Three-Dimensional Structures to
634 Protein Networks Provides Evolutionary Insights. Science. 314:1938–1941. doi:
635 10.1126/science.1136174.
- 636 Kim R, Guo J. 2010. Systematic analysis of short internal indels and their impact on
637 protein folding. BMC Struct Biol. 10:24. doi: 10.1186/1472-6807-10-24.
- 638 Kimura M. 1983. *The neutral theory of molecular evolution*. Cambridge University Press:
639 Cambridge.
- 640 Lek M et al. 2016. Analysis of protein-coding genetic variation in 60,706 humans.
641 Nature. 536:285–291. doi: 10.1038/nature19057.
- 642 Levy Karin E, Susko E, Pupko T. 2014. Alignment Errors Strongly Impact Likelihood-
643 Based Tests for Comparing Topologies. Molecular Biology and Evolution.
644 31:3057–3067.
- 645 Light S, Sagit R, Ekman D, Elofsson A. 2013. Long indels are disordered: A study of
646 disorder and indels in homologous eukaryotic proteins. Biochimica et Biophysica
647 Acta (BBA) - Proteins and Proteomics. 1834:890–897. doi:
648 10.1016/j.bbapap.2013.01.002.
- 649 Light S, Sagit R, Sachenkova O, Ekman D, Elofsson A. 2013. Protein Expansion Is
650 Primarily due to Indels in Intrinsically Disordered Regions. Molecular Biology and
651 Evolution. 30:2645–2653. doi: 10.1093/molbev/mst157.
- 652 Lin M et al. 2017. Effects of short indels on protein structure and function in human
653 genomes. Sci Rep. 7:9313. doi: 10.1038/s41598-017-09287-x.
- 654 Liu S et al. 2016. A facile and efficient transposon mutagenesis method for generation
655 of multi-codon deletions in protein sequences. Journal of Biotechnology. 227:27–
656 34. doi: 10.1016/j.jbiotec.2016.03.038.
- 657 Liu S et al. 2015. Structural plasticity of green fluorescent protein to amino acid
658 deletions and fluorescence rescue by folding-enhancing mutations. BMC
659 Biochemistry. 16:17. doi: 10.1186/s12858-015-0046-5.

660 Lynch M. 2007. *The origins of genome architecture*.
661 <https://repository.library.georgetown.edu/handle/10822/548280> (Accessed
662 August 11, 2023).

663 Mariani V, Biasini M, Barbato A, Schwede T. 2013. IDDT: a local superposition-free
664 score for comparing protein structures and models using distance difference
665 tests. *Bioinformatics*. 29:2722–2728. doi: 10.1093/bioinformatics/btt473.

666 Marsh JA, Teichmann SA. 2014. Parallel dynamics and evolution: Protein
667 conformational fluctuations and assembly reflect evolutionary changes in
668 sequence and structure. *BioEssays*. 36:209–218. doi: 10.1002/bies.201300134.

669 Mi H et al. 2017. PANTHER version 11: expanded annotation data from Gene Ontology
670 and Reactome pathways, and data analysis tool enhancements. *Nucleic acids
671 research*. 45:D183–D189. doi: 10.1093/nar/gkw1138.

672 Mi H, Muruganujan A, Ebert D, Huang X, Thomas PD. 2019. PANTHER version 14:
673 more genomes, a new PANTHER GO-slim and improvements in enrichment
674 analysis tools. *Nucleic acids research*. 47:D419–D426. doi:
675 10.1093/nar/gky1038.

676 Mi H, Muruganujan A, Thomas PD. 2013. PANTHER in 2013: modeling the evolution of
677 gene function, and other gene attributes, in the context of phylogenetic trees.
678 *Nucleic Acids Res*. 41:D377–386. doi: 10.1093/nar/gks1118.

679 Mills RE et al. 2011. Natural genetic variation caused by small insertions and deletions
680 in the human genome. *Genome Res*. 21:830–839. doi: 10.1101/gr.115907.110.

681 Mintseris J, Weng Z. 2005. Structure, function, and evolution of transient and obligate
682 protein–protein interactions. *Proceedings of the National Academy of Sciences*.
683 102:10930–10935. doi: 10.1073/pnas.0502667102.

684 Mitternacht S. 2016. FreeSASA: An open source C library for solvent accessible surface
685 area calculations. *F1000Res*. 5:189. doi: 10.12688/f1000research.7931.1.

686 Montgomery SB et al. 2013. The origin, evolution, and functional impact of short
687 insertion–deletion variants identified in 179 human genomes. *Genome Res*.
688 23:749–761. doi: 10.1101/gr.148718.112.

689 Nielsen R, Slatkin M. 2013. *An introduction to population genetics: theory and
690 applications*. Sinauer Associates Sunderland, MA.

691 Nishizawa M, Nishizawa K. 2002. A DNA Sequence Evolution Analysis Generalized by
692 Simulation and the Markov Chain Monte Carlo Method Implicates Strand
693 Slippage in a Majority of Insertions and Deletions. *J Mol Evol*. 55:706–717. doi:
694 10.1007/s00239-002-2366-5.

- 695 Ohta T. 1972. Evolutionary rate of cistrons and DNA divergence. *J Mol Evol.* 1:150–7.
696 doi: 10.1007/BF01659161.
- 697 Pascarella S, Argos P. 1992. Analysis of insertions/deletions in protein structures.
698 *Journal of Molecular Biology.* 224:461–471. doi: 10.1016/0022-2836(92)91008-D.
- 699 Penn O, Privman E, Ashkenazy H, et al. 2010. GUIDANCE: a web server for assessing
700 alignment confidence scores. *Nucleic acids research.* 38:W23–W28.
- 701 Penn O, Privman E, Landan G, Graur D, Pupko T. 2010. An alignment confidence score
702 capturing robustness to guide tree uncertainty. *Mol Biol Evol.* 27:1759–1767.
- 703 Privman E, Penn O, Pupko T. 2012. Improving the performance of positive selection
704 inference by filtering unreliable alignment regions. *Molecular Biology and*
705 *Evolution.* 29:1–5. doi: 10.1093/molbev/msr177.
- 706 Rockah-Shmuel L et al. 2013. Correlated Occurrence and Bypass of Frame-Shifting
707 Insertion-Deletions (InDels) to Give Functional Proteins. *PLOS Genetics.*
708 9:e1003882. doi: 10.1371/journal.pgen.1003882.
- 709 Salari R, Schönhuth A, Hormozdiari F, Cherkasov A, Sahinalp SC. 2008. The Relation
710 between Indel Length and Functional Divergence: A Formal Study. In: *Algorithms*
711 *in Bioinformatics.* Crandall, KA & Lagergren, J, editors. Lecture Notes in
712 *Computer Science* Springer: Berlin, Heidelberg pp. 330–341. doi: 10.1007/978-3-
713 540-87361-7_28.
- 714 Savino S, Desmet T, Franceus J. 2022. Insertions and deletions in protein evolution and
715 engineering. *Biotechnology Advances.* 60:108010. doi:
716 10.1016/j.biotechadv.2022.108010.
- 717 Scherrer MP, Meyer AG, Wilke CO. 2012. Modeling coding-sequence evolution within
718 the context of residue solvent accessibility. *BMC Evolutionary Biology.* 12:179.
719 doi: 10.1186/1471-2148-12-179.
- 720 Shahmoradi A et al. 2014. Predicting Evolutionary Site Variability from Structure in Viral
721 Proteins: Buriedness, Packing, Flexibility, and Design. *J Mol Evol.* 79:130–142.
722 doi: 10.1007/s00239-014-9644-x.
- 723 Shih C-H, Chang C-M, Lin Y-S, Lo W-C, Hwang J-K. 2012. Evolutionary information
724 hidden in a single protein structure. *Proteins: Structure, Function, and*
725 *Bioinformatics.* 80:1647–1657. doi: 10.1002/prot.24058.
- 726 Simm AM, Baldwin AJ, Busse K, Jones DD. 2007. Investigating protein structural
727 plasticity by surveying the consequence of an amino acid deletion from TEM-1 β -
728 lactamase. *FEBS Letters.* 581:3904–3908. doi: 10.1016/j.febslet.2007.07.018.

729 Sinden RR et al. 2002. Triplet repeat DNA structures and human genetic disease:
730 dynamic mutations from dynamic DNA. *J Biosci.* 27:53–65. doi:
731 10.1007/BF02703683.

732 Slodkowicz G, Goldman N. 2020. Integrated structural and evolutionary analysis reveals
733 common mechanisms underlying adaptive evolution in mammals. *Proc Natl Acad*
734 *Sci USA.* 117:5977–5986. doi: 10.1073/pnas.1916786117.

735 Snir S, Pachter L. 2006. Phylogenetic Profiling of Insertions and Deletions in Vertebrate
736 Genomes. In: *Research in Computational Molecular Biology.* Apostolico, A,
737 Guerra, C, Istrail, S, Pevzner, PA, & Waterman, M, editors. *Lecture Notes in*
738 *Computer Science* Springer: Berlin, Heidelberg pp. 265–280. doi:
739 10.1007/11732990_23.

740 Tao S et al. 2007. Patterns of Insertion and Deletion in Mammalian Genomes. *Current*
741 *Genomics.* 8:370–378. doi: 10.2174/138920207783406479.

742 Taylor MS, Ponting CP, Copley RR. 2004. Occurrence and Consequences of Coding
743 Sequence Insertions and Deletions in Mammalian Genomes. *Genome Res.*
744 14:555–566. doi: 10.1101/gr.1977804.

745 The Genomes Project C. 2015. A global reference for human genetic variation. *Nature.*
746 526:68–74. doi: 10.1038/nature15393
747 [http://www.nature.com/nature/journal/v526/n7571/abs/nature15393.html#supple](http://www.nature.com/nature/journal/v526/n7571/abs/nature15393.html#supplementary-information)
748 [mentary-information.](http://www.nature.com/nature/journal/v526/n7571/abs/nature15393.html#supplementary-information)

749 Thomas PD et al. 2022. PANTHER: Making genome-scale phylogenetics accessible to
750 all. *Protein Science.* 31:8–22. doi: 10.1002/pro.4218.

751 Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. 2013. Maximum Allowed
752 Solvent Accessibilities of Residues in Proteins. *PLoS One.* 8:e80635. doi:
753 10.1371/journal.pone.0080635.

754 Tóth-Petróczy Á, Tawfik DS. 2014. Hopeful (Protein InDel) Monsters? *Structure.*
755 22:803–804. doi: 10.1016/j.str.2014.05.013.

756 Tóth-Petróczy Á, Tawfik DS. 2011. Slow protein evolutionary rates are dictated by
757 surface–core association. *Proceedings of the National Academy of Sciences.*
758 108:11151–11156. doi: 10.1073/pnas.1015994108.

759 Varadi M et al. 2022. AlphaFold Protein Structure Database: massively expanding the
760 structural coverage of protein–sequence space with high-accuracy models.
761 *Nucleic acids research.* 50:D439–D444.

762 Won YJ, Hey J. 2005. Divergence population genetics of chimpanzees. *Mol Biol Evol.*
763 22:297–307. doi: 10.1093/molbev/msi017.

- 764 Woods H et al. 2023. Computational modeling and prediction of deletion mutants.
765 Structure. 31:713-723.e3. doi: 10.1016/j.str.2023.04.005.
- 766 Yeh S-W et al. 2014. Site-Specific Structural Constraints on Protein Sequence
767 Evolutionary Divergence: Local Packing Density versus Solvent Exposure.
768 Molecular Biology and Evolution. 31:135–139. doi: 10.1093/molbev/mst178.
- 769 Zhang Z, Huang J, Wang Z, Wang L, Gao P. 2011. Impact of Indels on the Flanking
770 Regions in Structural Domains. Molecular Biology and Evolution. 28:291–301.
771 doi: 10.1093/molbev/msq196.
- 772 Zhang Z, Wang J, Gong Y, Li Y. 2018. Contributions of substitutions and indels to the
773 structural variations in ancient protein superfamilies. BMC Genomics. 19:771.
774 doi: 10.1186/s12864-018-5178-8.
- 775 Zhang Z, Wang Y, Wang L, Gao P. 2010. The Combined Effects of Amino Acid
776 Substitutions and Indels on the Evolution of Structure within Protein Families.
777 PLOS ONE. 5:e14316. doi: 10.1371/journal.pone.0014316.
- 778 Zhao Z et al. 2000. Worldwide DNA sequence variation in a 10-kilobase noncoding
779 region on human chromosome 22. Proc Natl Acad Sci USA. 97:11354–8. doi:
780 10.1073/pnas.200348197.

781
782
783

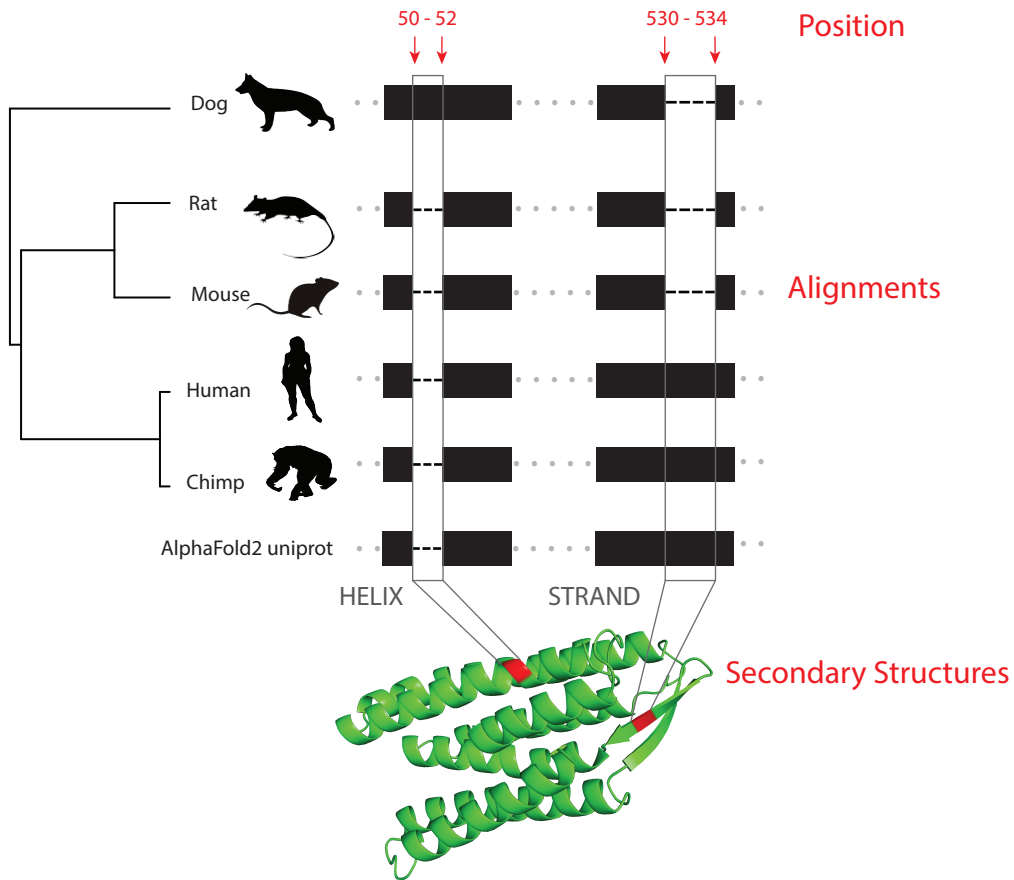


Figure 2

ALL
11444 genes, 97383 indels, 200 iterations

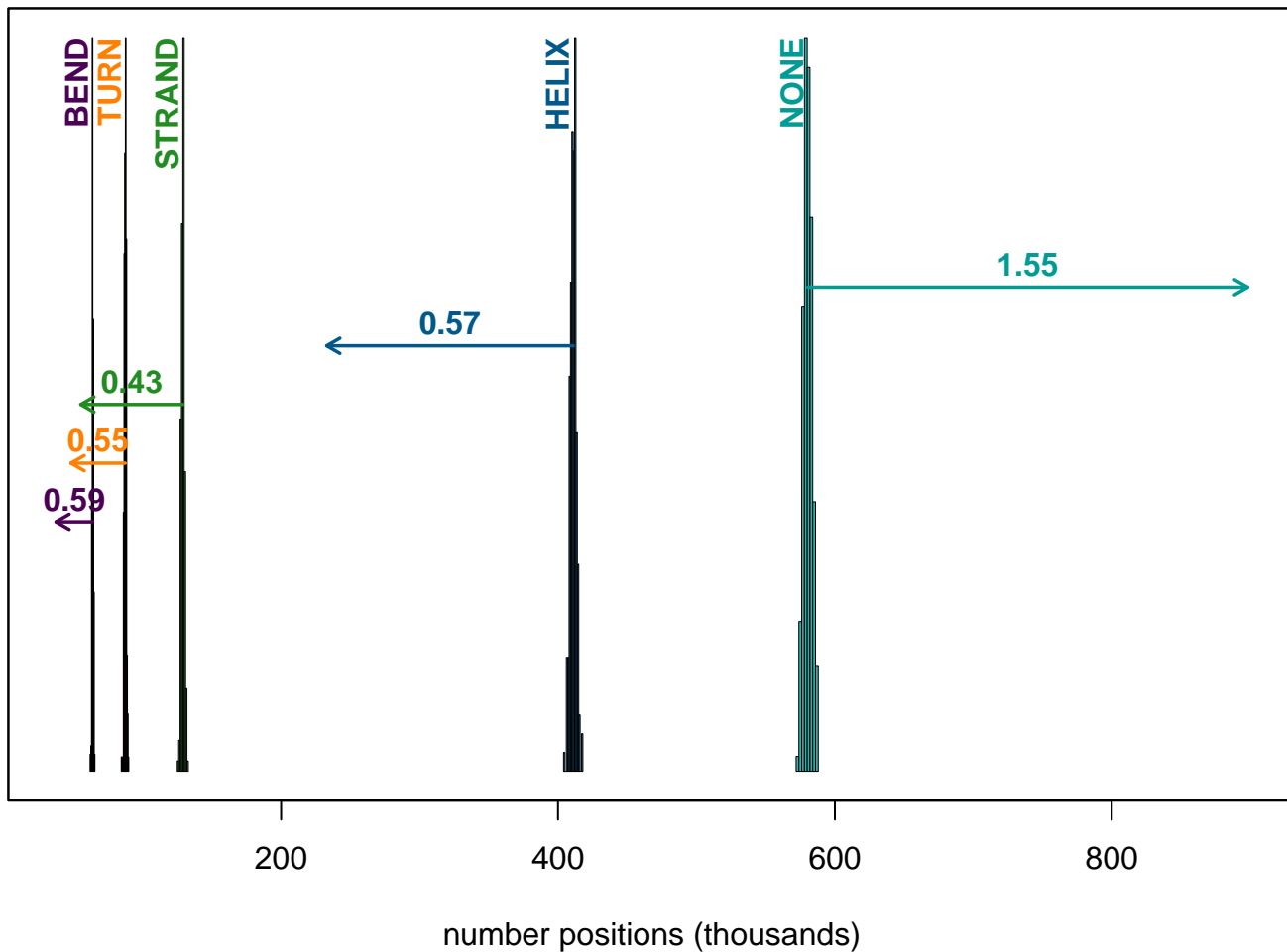


Figure 3

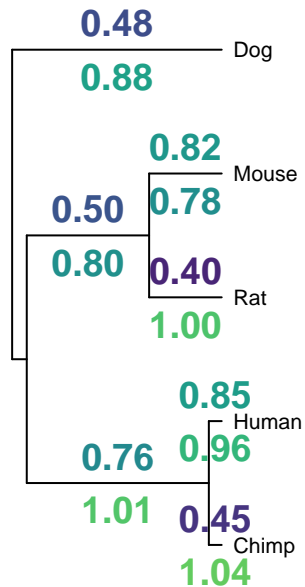
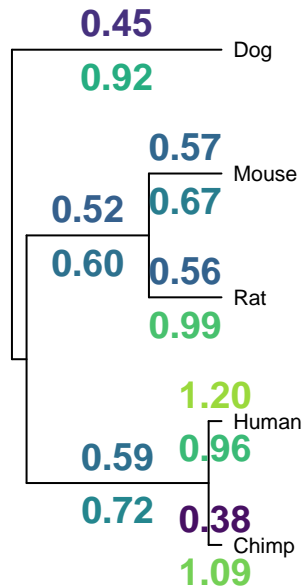
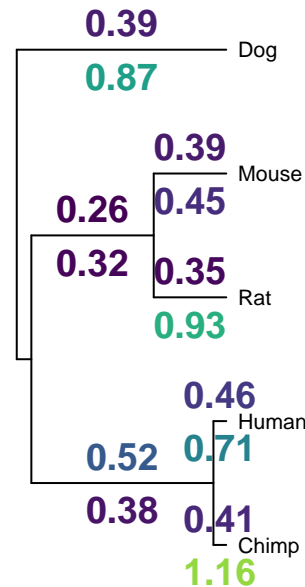
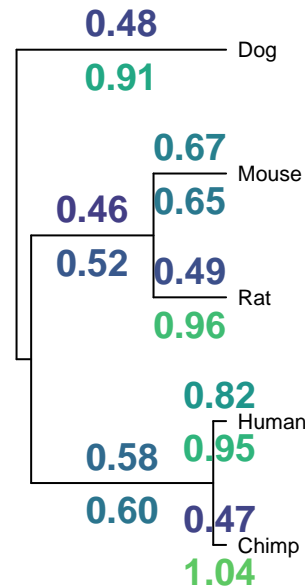
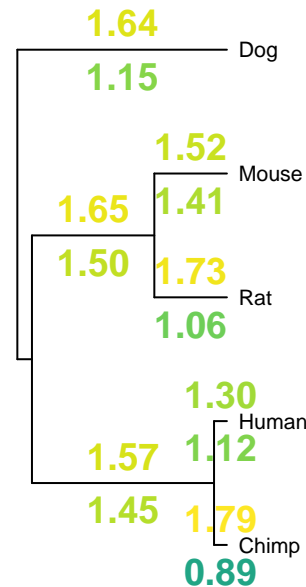
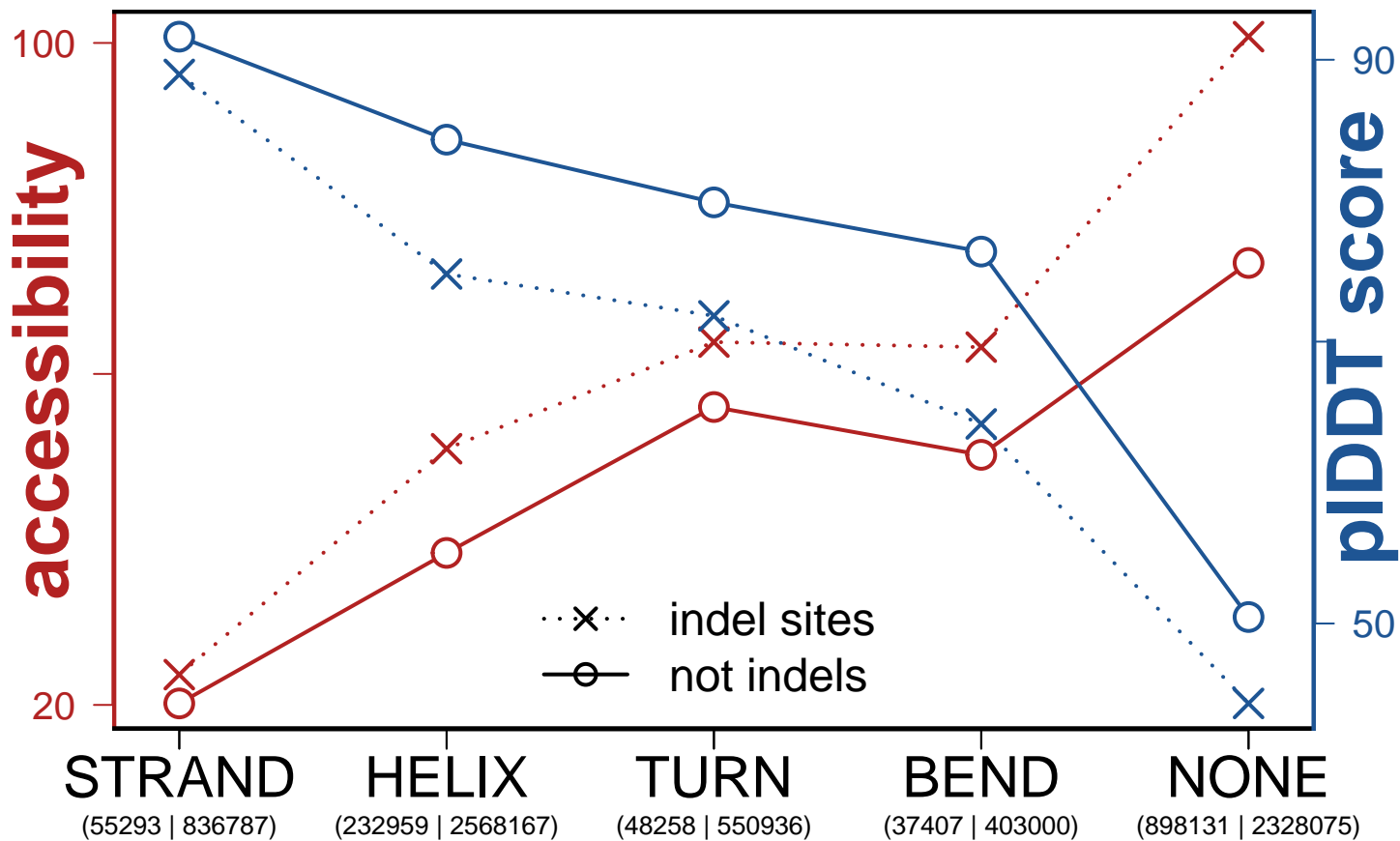
BEND**TURN****STRAND****HELIX****NONE****ALL**

Figure 4

ALL



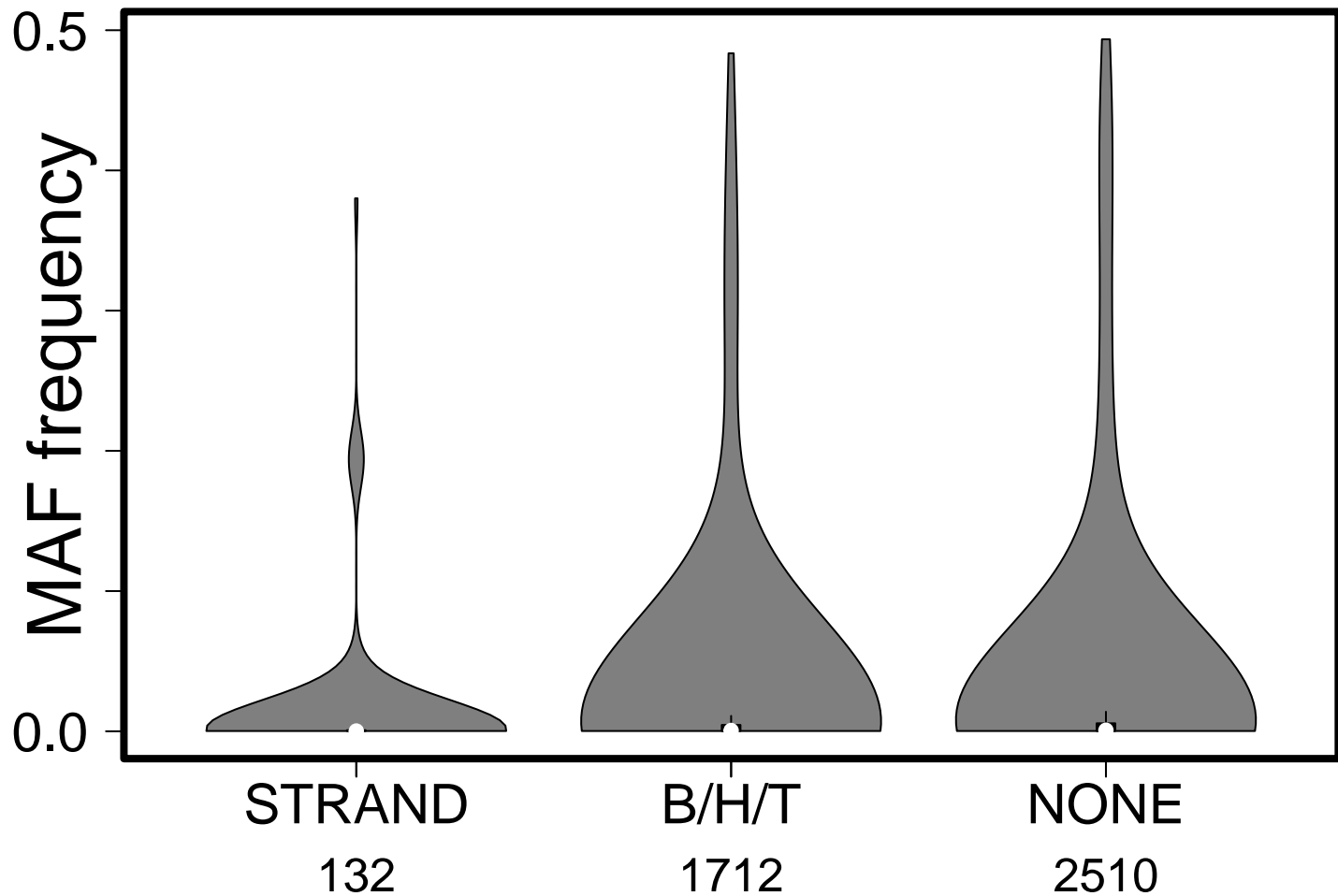


Table 1. Number of indels or codon mutations that overlap secondary structures. Observed=number of positions in alignments that map over each category. Expected=Number expected based on randomization. Codons are classified as invariant (Invariant), synonymous (Syn.) or nonsynonymous (Non.). p=proportion of sites within their respective columns that fall within each category. This table is repeated as Supplementary Tables 2, after employing four different subsetting strategies.

	Indels			Codon-based							
	Observed	Expected	O/E	Invariant	p Inv.	Syn.	p Syn.	Syn./Inv.	Non.	p Non.	Non./Inv.
STRAND	55,293	129,070	0.43	455,059	0.120	278,936	0.130	1.09	143,454	0.086	0.72
TURN	48,258	87,473	0.55	287,034	0.076	189,311	0.088	1.17	110,149	0.066	0.87
HELIX	232,959	411,110	0.57	1,381,189	0.364	827,926	0.386	1.06	532,917	0.320	0.88
BEND	37,407	63,890	0.59	209,328	0.055	137,150	0.064	1.16	84,632	0.051	0.92
NONE	898,131	580,490	1.55	1,464,044	0.386	709,265	0.331	0.86	796,815	0.478	1.24

Table 2. AUC metrics for three Generalized Linear Models. Mean (standard deviation) AUC from 5 iterations of randomly sampling sites across alignments.

analysis_type	indel~SS	indel~RSA	indel~SS+RSA
ALL	0.684 (0.004)	0.707 (0.008)	0.720 (0.009)
INTERNAL	0.614 (0.010)	0.604 (0.007)	0.612 (0.005)
GU94_PA100_GD40	0.622 (0.006)	0.597 (0.015)	0.610 (0.013)
LENGTH_LTE20	0.618 (0.009)	0.610 (0.010)	0.621 (0.011)
MERGED	0.618 (0.006)	0.610 (0.014)	0.618 (0.011)