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**

A fundamental goal in evolutionary biology and population genetics is to understand

how selection shapes the fate of new mutations. Here we test the null hypothesis that

insertion-deletion events (indels) in protein coding regions occur randomly with respect

to secondary structures. We identified indels across 11,444 sequence alignments in

mouse, rat, human, chimp, and dog genomes, then quantified their overlap with four

different types of secondary structure – alpha helices, beta strands, protein bends, and

protein turns – predicted by deep-learning methods of AlphaFold2. Indels overlapped

secondary structures 54% as much as expected, and were especially under-

represented over beta strands, which tend to form internal, stable regions of proteins. In

contrast, indels were enriched by 155% over regions without any predicted secondary

structures. These skews were stronger in the rodent lineages compared to the primate

lineages, consistent with population genetic theory predicting that natural selection will

be more efficient in species with larger effective population sizes. Nonsynonymous

substitutions were also less common in regions of protein secondary structure, although

not as strongly reduced as in indels. In a complementary analysis of thousands of

human genomes, we showed that indels overlapping secondary structure segregated at

significantly lower frequency than indels outside of secondary structure. Taken together,

our study shows that indels are selected against if they overlap secondary structure,

presumably because they disrupt the tertiary structure and function of a protein.

36

## Significance

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

44

## Introduction

Understanding the fate of new mutations is critical to defining the evolutionary processes that shape biological diversity. At the level of single nucleotides, a rich body of theory has been developed to infer whether mutations are neutral, deleterious, or beneficial (reviewed by Hedrick 2005; Hartl and Clark 2007; Nielsen and Slatkin 2013). Understanding the selective impact of insertion-deletion events (indels), which can extend many nucleotides, has proven to be much more complicated.

Previous studies investigating the functional impact of indels generally fall into two categories (Savino et al. 2022). First, protein engineering studies have shown that indels can impact a protein's function, especially if they overlap important secondary structures (Simm et al. 2007; Arpino et al. 2014; Tóth-Petróczy and Tawfik 2014; Gavrilov et al. 2015; Grocholski et al. 2015; Liu et al. 2015; Liu et al. 2016; Jackson et al. 2017; Gavrilov et al. 2018; Halliwell et al. 2018; Gonzalez et al. 2019; Woods et al. 2023). For example, Liu et al. (2016) found that experimentally deleting amino acids in beta strands and alpha helices of Green Fluorescent Protein tended to reduce fluorescence, while deletions outside such regions were relatively neutral.

Second, evolutionary and population genetic studies have suggested that indels are relatively deleterious if they are long (Pascarella and Argos 1992; Taylor et al. 2004; Tao et al. 2007; Hsing and Cherkasov 2008; Kim and Guo 2010; Mills et al. 2011; Rockah-Shmuel et al. 2013; Lek et al. 2016; Zhang et al. 2018), cause frame-shifts (Iengar 2012; Chong et al. 2013; Montgomery et al. 2013; Bermejo-Das-Neves et al. 2014; Chen and Guo 2021), occur internally in the protein (Lin et al. 2017), alter flanking amino acids (Zhang et al. 2011), or fall outside of disordered regions (Taylor et al. 2004;

68  Light, Sagit, Ekman, et al. 2013; Light, Sagit, Sachenkova, et al. 2013; Bermejo-Das-

69  Neves et al. 2014; Khan et al. 2015). Protein families with indels tend to diverge in their

70  structure and function relative to protein families without indels (Salari et al. 2008;

71  Hormozdiari et al. 2009; Zhang et al. 2010; Gavrilov et al. 2015; Gavrilov et al. 2018;

72  Zhang et al. 2018; Banerjee et al. 2019; Jayaraman et al. 2022), suggesting indels can

73  be an important source of evolutionary novelty. Indeed, one study estimated that >70%

74  of indels that have reached fixation have done so through positive selection (Barton and

75  Zeng 2019).

76      Two important evolutionary studies identified orthologs across species and then

77  overlapped inferred indels with experimentally determined protein structures in the

78  Protein Data Bank (PDB, Berman et al. 2000). Following the publication of the human,

79  mouse and rat genomes, Taylor et al. (2004) identified 52 orthologous protein-coding

80  genes that had an indel *and* a protein structure. Of these 52 indels, 31.5% of their

81  sequence overlapped secondary structure of any kind, compared to 52.5% expected. A

82  few years later, de la Chaux et al. (2007) analyzed the distribution of 343 protein-coding

83  indels identified from human-chimp-rhesus orthologs that also occurred in the PDB.

84  They found a deficiency of indels that overlapped alpha helices, but no difference in

85  indels that overlapped beta strands.

86      As impactful as these studies were, they may not paint a full picture of the

87  functional consequences of indel variation. The set of genes that could be studied was

88  small, mostly limited by structural protein data or annotated Pfam domains. Pfam

89  domains do not necessarily correlate with 3D structure and the PDB represents a

90  biased set of proteins (or protein regions) that are amenable to the experimental

91  approaches required for structural proteomics, such as their ability to be crystallized.

92  The relatively biased set of proteins for which we have structural data thus limits a

93  systematic analysis across full genomes. For example, one study of duplicated genes

94  could not analyze full-length proteins because of divergence between aligned gene

95  sequences and proteins represented in the PDB (Guo et al. 2012). However, the recent

96  release of AlphaFold2 – a deep-learning project that accurately predicts the 3D

97  structure of a protein from its amino acid sequence (Jumper et al. 2021; Varadi et al.

98  2022) – provides a unique opportunity to systematically study indels across full proteins

99  and whole genomes.

100  Here we combine genome-wide predictions of AlphaFold2 with evolutionary and

101  population genetic methods to ask whether indels occur randomly with respect to

102  secondary structure, providing the most comprehensive evolutionary investigation into

103  the fate of indels in protein coding regions. We report four main results: 1) 97,382 indels

104  identified from 11,444 five-species alignments in the tree (dog, ((mouse, rat), (human,

105  chimp)) overlapped secondary structures 54% as often as expected, but were 155%

106  more common than expected in regions with no predicted secondary structures, 2)

107  indels that overlapped beta strands and occurred internally in a protein were especially

108  rare, consistent with the known importance of these regions in overall protein structure,

109  3) skews in observed vs. expected were stronger in the rodent lineages compared to

110  the primate lineages, consistent with theory predicting more efficient selection in rodents

111  given their larger effective population sizes, and 4) within human populations, indels that

112  overlapped secondary structures occurred at significantly lower frequency compared to

113  indels outside of secondary structures. Taken together, our results indicate selection

114     acts against indels when they arise over structurally important regions of proteins,

115     presumably because they can disrupt overall structure and therefore the function of a

116     protein.

117

118                         **Materials and Methods**

119     **Interspecific insertion-deletion (indel) events.** We downloaded protein sequences

120     from all protein-coding genes identified as one-to-one orthologs between mouse, rat,

121     human, chimp, and dog from Ensembl version 107 (ensembl.org). In the case of

122     alternative transcripts, we chose the longest translated transcript to represent the gene.

123     11,444 genes had one-to-one orthologs across all five species.

124          We aligned proteins using GUIDANCE (Penn, Privman, Landan, et al. 2010; Penn,

125     Privman, Ashkenazy, et al. 2010; Privman et al. 2012; Levy Karin et al. 2014). This

126     approach estimates per-site alignment confidence by calculating its consistency across

127     different starting guide trees, allowing us to incorporate a measure of confidence in

128     downstream analyses. Importantly, we could use GUIDANCE scores to estimate error in

129     indel placement and identify indels that were confidently placed. In each GUIDANCE

130     iteration, we aligned protein sequences with MAFFT (Katoh et al. 2002). We ran MAFFT

131     under the recommended default parameters; in the case of indels the most important

132     default parameters were the gap opening penalty (default=1.53) and gap offset value

133     (similar to gap extension penalty, default=0.123). We then identified all indels as gaps

134     from all 11,444 alignments (Fig. 1).

135          Our analyses could be impacted by sequencing errors or annotation errors that

136     result in spurious inclusion or exclusion of amino acids from certain genes, or by

137    alignment errors (Fitch and Smith 1983; Chowdhury and Garai 2017). Therefore, we

138    repeated all downstream analyses after subsetting indels in four different ways: 1)

139    INTERNAL: any indels that reached the beginning or ends of alignments were excluded,

140    as visual inspection indicated these were noisy regions of alignment that could be

141    related to incomplete annotation of full length genes, 2) GU94_PA100_GD40:

142    INTERNAL indels whose flanking five positions on both 5' and 3' ends (10 flanking

143    positions total) had an average GUIDANCE confidence score of at least 0.94 (median

144    observed), contained no overlapping indels, and had an average Grantham distance

145    (Grantham 1974) of less than 40 (median observed), where Grantham distance was

146    calculated using the R package AGVGD (https://CRAN.R-project.org/package=agvgd).

147    This subset was meant to enrich for well-anchored indels and avoid problems

148    distinguishing gaps in alignment due to protein divergence, versus gaps in alignment to

149    insertion-deletion events (Snir and Pachter 2006; Salari et al. 2008; Jilani et al. 2022),

150    3) LENGTH_LTE20: INTERNAL indels that were less than or equal to 20 amino acids

151    long in length, minimizing the impact of large indels that sometimes appeared to be

152    spurious, and 4) MERGED: INTERNAL indels after merging coordinates that

153    overlapped, so that sites in an alignment that were in different overlapping regions only

154    contributed once. We present the results from these four subsets as supplementary

155    files, but they all produced essentially identical results as analyzing ALL indels.

156

157    **AlphaFold2.** AlphaFold2 is a deep learning approach developed by DeepMind to

158    predict the 3D structure of proteins from only their amino acid sequence (Jumper et al.

159  2021; Varadi et al. 2022). Comparison to empirical data indicates these computational

160  predictions are over 90% accurate.

161      AlphaFold2 assigns 43 different secondary structures to different regions of a

162  protein, which we collapsed into five main categories. There were 32 different

163  AlphaFold2 predictions that contained the phrase HELX, which are predictions of

164  different helices; we collapsed these into the single term HELIX. There were 8 different

165  AlphaFold2 predictions that contain the phrase TURN, which are regions where the

166  polypeptide is predicted to reverse direction in 3D space; we collapsed these into the

167  single term TURN. We included the single Alphafold2 prediction STRAND as-is, which

168  are regions predicted to contain beta strands (also referred to as beta sheets). We

169  included the single AlphaFold2 prediction BEND as-is, which are regions where the

170  polypeptide is predicted to change direction but not fully reverse. There was one last

171  Alphafold2 prediction OTHER, but we did not observe any instances of this prediction in

172  any of the proteins analyzed in this study so ignored that term. Each residue in the

173  Uniprot protein used by AlphaFold2 was assigned to one of these four categories, or

174  assigned the term NONE if they occurred outside any predicted secondary structure.

175      To link AlphaFold2 predictions to our five-species alignments above, we included

176  the Uniprot sequence in the alignment (Fig. 1). In rare cases, the AlphaFold2-

177  downloaded Uniprot sequence did not match the Ensembl-downloaded Uniprot

178  sequence, in which case we discarded the alignment from all analyses. Each position in

179  each indel was then assigned HELIX, STRAND, TURN, BEND, or NONE (Fig. 1). In

180  cases where the Uniprot sequence was "deleted" (for example, indel 50-52 in Fig. 1),

181     we assigned one-half of the deleted positions to whatever was assigned to its 5'-flanking

182     residue, and one-half to whatever was assigned to its 3'-flanking residue.

183

184     **Randomization of indel positions.** We generated null expectations through a

185     randomization procedure. For each alignment, we randomly shuffled the starting

186     position of each indel, then extended each randomized indel by its observed length. In

187     cases where a randomized indel extended past the end of an alignment, we wrapped

188     the randomized indel to the front of the alignment. After shuffling the unique indels

189     within each alignment, we re-calculated the number of residues falling in each

190     secondary structure, exactly as described above. We repeated this process 200 times

191     to generate null expectations.  We repeated this entire process for the four different

192     subsets described above. For these four subsets, the relevant alignments were first

193     truncated to match included regions and provide a more appropriate background for

194     randomization.

195

196     **Gene Ontology enrichment.** For the MERGED indels only, we identified relative

197     outliers by counting the number of sites in the alignment overlapping NONE vs. not,

198     versus sites overlapping indels vs. not. We excluded alignments that had fewer than 5

199     positions in any of these four cells of this 2x2 table, then applied a $X^2$ test and corrected

200     resulting p-values (Benjamini and Hochberg 1995). Genes with a -log10 p.value of at

201     least 10 and at least a 1.5 fold change in expectation were taken as relative outliers. We

202     tested whether these relative outlier genes were enriched for any Biological Process,

203     Molecular Function, or Cellular Component using Panther Classification system (Mi et

204 al. 2013; Mi et al. 2017; Mi et al. 2019; Thomas et al. 2022), run from PantherDB

205 (https://pantherdb.org/), with the settings "Test Type=Fisher's Exact Test" and

206 "Correction=Calculate False Discovery Rate". We also performed Gene Ontology

207 analyses for genes which had no indels across the five species analyzed.

208

209 **Accessibility and pIDDT scores.** Sites that are relatively internal on a 3D

210 protein evolve more slowly than external sites, both at the level of nonsynonymous

211 mutations (Goldman et al. 1998; Bustamante et al. 2000; Dean et al. 2002; Franzosa

212 and Xia 2009; Tóth-Petróczy and Tawfik 2011; Scherrer et al. 2012; Shih et al. 2012;

213 Marsh and Teichmann 2014; Shahmoradi et al. 2014; Yeh et al. 2014) and indel

214 variation (Hsing and Cherkasov 2008; Guo et al. 2012). This correlation is complicated

215 by whether or not external residues interact with other proteins (Mintseris and Weng

216 2005; Kim et al. 2006), or if externally oriented residues form active sites of proteins

217 (Slodkowicz and Goldman 2020). For each site in each alignment, we calculated

218 relative solvent accessibility, which is the degree to which a residue occurs on the

219 outside of a folded protein (Tien et al. 2013), using FREESASA (Mitternacht 2016) with

220 the "--format=rsa" option, using the AlphaFold2 structure as input. We also compared

221 pIDDT scores (Mariani et al. 2013) across an alignment. pIDDT scores are

222 computational measures of confidence included in AlphaFold2 predictions. According to

223 AlphaFold2, pIDDT scores <50 likely represent intrinsically disordered or unstructured

224 regions. As above, any "deletions" in the Uniprot sequence were divided, and one-half

225 of their sites were assigned the accessibility and pIDDT scores of their 5' flanking

226 residue, and the other half to the scores of their 3' flanking residue.

227    As will be shown below, secondary structure and relative solvent accessibility are

228    strongly correlated. In an attempt to separate the effects of these two features on the

229    probability of observing an indel, we compared Receiver Operating Characteristic

230    (ROC) curves and Area Under the Curve (AUC) values from three Generalized Linear

231    Models and then compared their likelihoods. Two models tested whether the probability

232    of observing an indel was a function of secondary structure or relative solvent

233    accessibility alone – glm(indel~secondary_structure) or glm(indel~rsa), respectively. A

234    third model included both as independent variables – glm(indel~secondary_structure +

235    rsa). We quantified the gain in likelihood when we included both independent variables,

236    versus each one separately. For all three models we included the "family = binomial"

237    argument to model logistic variance. Our approach closely followed that of Jackson et

238    al. (2017), modifying their scripts to suit our approach.

239    Because sites in a protein are not independent from each other, before applying

240    Generalized Linear Models we randomly sampled a single site from each alignment.

241    However, we did not sample sites with equal probability. Instead, we downweighted the

242    probability of sampling by the inverse of the grand total of the five secondary structures

243    (HELIX, STRAND, TURN, BEND, or NONE). By including this weighting scheme, we

244    ensured even sampling of secondary structures, increasing power of all three

245    Generalized Linear Models.

246

247    **Comparison to synonymous and nonsynonymous mutations.** To provide additional

248    context with which to interpret the distribution of indels, we tested three different

249    nucleotide-based sites. First, we quantified the distribution of invariant sites across

250   secondary structure as a kind of null distribution. Then we quantified the same with

251   respect to synonymous and nonsynonymous sites. We predicted that synonymous sites

252   should distribute similarly to invariant sites, because they do not alter the protein

253   sequence and thus probably have relatively minor effect on secondary structure.

254   Conversely, we predicted that nonsynonymous sites would occur less frequently over

255   secondary structure because, all else equal, their resulting amino acid changes could

256   alter secondary structure.

257        Using the same 5-species alignments above, we reverse-translated each protein

258   to its transcript, downloaded from Ensembl version 107. We counted the proportion of

259   synonymous vs. nonsynonymous variants occurring over the different secondary

260   structures, compared to invariant sites. We only quantified synonymous vs.

261   nonsynonymous variants from the same alignments and sites that were used in our

262   indel analyses.

263

264   **Intraspecific indel events.** As a complementary analysis to the interspecific analyses

265   described above, we analyzed intraspecific variation from Phase 3 of the 1000 Human

266   Genomes project (https://www.internationalgenome.org/data-portal/data-collection/30x-

267   grch38) (The Genomes Project 2015; Byrska-Bishop et al. 2022). This database

268   contains haplotype-phased indel calls (files named like

269   ALL.chr1.shapeit2_integrated_snvindels_v2a_27022019.GRCh38.phased.INDELS.vcf)

270   from 2,504 unrelated samples from 26 populations, with sample size ranging from 61 to

271   113 per population. These 26 populations derive from five large geographic areas:

272   Africa, East Asia, South Asia, South America, and Europe.

273        Indel coordinates were truncated to match exon coordinates downloaded from

274    UCSC Table Browser (table name=unipAliSwissprot from GRCH38). For any protein-

275    coding genes that contained at least one indel, we assembled the reference and

276    alternative alleles from the human genome, computationally placed indels, and then

277    translated both alleles. Any indels that resulted in a frameshift in the first 95% of the

278    protein-coding transcript (counted from 5' translation start site) were excluded, because

279    it is unclear whether reference and alternative alleles share 3D structure if they are

280    dramatically frame-shifted with respect to each other.

281        We only analyzed genes that were part of the five-species interspecific analyses

282    described above. Otherwise, we would have included recent human-specific duplicates,

283    where predictions might become noisy because of uncertainty about the exact timing of

284    duplication along the lineage to modern humans.

285

286                                           **Results**

287    **Indels were depleted in regions with secondary structure.** There were 11,444

288    genes that had one-to-one orthologs between dog, mouse, rat, chimp, and human

289    genomes. Across these 11,444 alignments we identified 97,382 indels spanning

290    1,272,048 positions. Indel sizes ranged from 1 to 2,870 residues long, but most were

291    small: the 25%, 50%, and 75% quantiles were 1, 3, and 10 residues, respectively. Indel

292    positions overlapped secondary structures significantly less than expected (Fig. 2, Table

293    1). Indel positions were most under-represented in STRAND, occurring at 43%

294    expectations (calculated as 55,293 indel sites that overlapped STRAND, compared to

295    129,070 averaged across 200 randomizations), followed by indel positions occurring in

296     TURN (55%), HELIX (57%), and BEND (59%) (Table 1). In contrast, indel positions

297     occurred at 155% expectation in NONE, meaning indels were much more likely occur in

298     protein regions with no predicted secondary structure (Table 1). All observed values fell

299     far outside the distributions from randomization (Fig. 2), translating into a p-value of

300     essentially 0. We reached nearly identical conclusions after subsetting indels in four

301     ways described above (Supplementary Figure 1, Supplementary Table 1), with one

302     exception: indels over TURN and BEND are not under-represented in the very stringent

303     subset GU94_PA100_GD40 (Supplementary Figure 1, Supplementary Table 1).

304

305     **Skews in indel distribution were stronger in rodents.** By using dog as an outgroup,

306     we polarized all indels into either an insertion or deletion and placed each indel event on

307     a specific branch in the phylogenetic tree, using simple parsimony. In other words, if

308     amino acid sequences existed for mouse and rat, but not for the other species, that

309     indel was mapped as an insertion on the branch leading to rodents.

310         There are seven branches on the phylogenetic tree analyzed here. Across the

311     four secondary structures (BEND, TURN, STRAND, and HELIX), 24 of 28 O:E values

312     were lower for insertions compared to deletions (Figure 3). Conversely, across NONE

313     sites all branches showed higher O:E for insertions compared to deletions. Taken

314     together, these results suggest that insertions over secondary structure are more

315     deleterious than deletions.

316         For the four secondary structures, O:E values were consistently lower in rodent

317     lineages compared to primate lineages. There are four secondary structure that can be

318     mapped to three rodent branches and three primate branches, where each branch

319     contains insertions and deletions, for a total of 48 O:E values in Figure 2. 46 of these 48

320     O:E values were lower in the rodent lineages compared to primate lineages. For

321     example, O:E values for insertions over STRAND in the three rodent lineages = 0.26,

322     0.39, and 0.35, while in primates the three values = 0.52, 0.46, and 0.41. Conversely,

323     O:E values for NONE sites tend to be higher in rodents compared to primates. In sum,

324     indels were especially unlikely to overlap secondary structures in rodents. All patterns

325     described held after analyzing the four different subsets of indels described above

326     (Supplementary Figure 2).

327

328     **GO analysis.** We identified 797 alignments (genes) where the enrichment of indels over

329     NONE was especially high. Compared to the rest of the 4,995 alignments, these 797

330     genes showed no statistical enrichment of Biological Process, but under the Cellular

331     Component and Molecular Function ontologies showed enrichment of terms associated

332     with cilia and ubiquitination. This enrichment lacks an obvious explanation.

333        We identified 88 alignments (genes) whose indels overlapped NONE much less

334     than expected. None of these 88 genes showed enrichment of Biological Process or

335     Molecular Function but showed enrichment of gene products localized to the nucleus

336     under Cellular Component. In sum, there were no striking or consistent patterns of

337     Gene Ontology enrichment associated with outlier genes in either direction.

338        We also analyzed the 904 genes which had no indels across any of the five

339     species in the alignment. GO analysis uncovered many functional terms associated with

340     neurotransmission, including synapse localization and synaptic transmission

341     (Supplementary Table 2). This result suggests that genes involved in neurotransmission

342    may be especially intolerant of indel mutations. Interestingly, genes involved in immune

343    response appeared to be under-represented among genes with no indels. This result

344    may indicate that immune genes undergo indel mutations more often than expected.

345

346    **Indels were enriched in regions with high accessibility and low pIDDT scores.**

347    Accessibility and pIDDT scores varied according to secondary structure. STRAND had

348    low accessibility and high pIDDT scores, indicating these secondary structures tend to

349    fall on the inside of proteins and are relatively stable (Fig. 4). On the other end of the

350    spectrum, NONE sites were much more accessible, with lower pIDDT scores, indicating

351    external and unstable regions of proteins (Fig. 4).

352         Importantly, sites that overlapped indels consistently showed higher accessibility

353    and lower pIDDT scores (compare X vs. O within each group, Fig. 4). In other words,

354    *within* each secondary structure, indels were more commonly observed at sites that

355    were relatively external and in relatively unstable regions, compared to sites that did not

356    overlap indels. Woods et al. (2023) found that experimentally deleting amino acids that

357    reside in regions of high pIDDT were most likely to have a deleterious effect on protein

358    function, providing an explanation for why we observe indels more frequently in regions

359    with low pIDDT scores. This pattern held across all four subsets of indels described

360    above (Supplementary Figure 3).

361         Comparing three different Generalized Linear Models demonstrated that the

362    effects of secondary structure were indistinguishable from the effects of relative solvent

363    accessibility (Table 2). In the ALL dataset, secondary structure performed about as well

364    as relative solvent accessibility (AUC=0.684 vs. 0.707, respectively), and including both

365  as independent variables had only minor improvement to AUC (0.720) compared to

366  single regressions. Similar results were obtained across the four subsets of data

367  described above (Table 2). This shows that secondary structure and relative solvent

368  accessibility are so correlated with each other that their effects cannot be meaningfully

369  separated.

370

371  **Nonsynonymous variants were also depleted in protein regions with secondary**

372  **structure.** Among the 11,444 alignments, we analyzed 3.8, 2.14, and 1.67 million

373  codons that were invariant, synonymous, or nonsynonymous, respectively (Table 1).

374  Synonymous codons overlapped secondary structures as often as invariant codons

375  (synonymous-to-invariant ratios ranging from 0.86 to 1.17, Table 1). In contrast,

376  nonsynonymous codons occurred far less frequently across the four secondary

377  structures (nonsynonymous-to-invariant ratios ranging from 0.71 to 0.92) and more over

378  NONE (nonsynonymous-to-invariant ratio of 1.24) (Table 1). These nonsynonymous-to-

379  invariant ratios were generally smaller in magnitude than the O:E ratios estimated from

380  indel distribution (Table 1). For example, indels occurred at 43% expectation over

381  STRAND, while nonsynonymous codons occurred at 71% "expectation" (Table 1).

382       Similar patterns emerged after analyzing the four subsets of indels

383  (Supplementary Table 1). The main exception was that nonsynonymous-to-invariant

384  ratios ranged from 0.91 to 0.98 across the four secondary structures, and from 1.05 to

385  1.09 for NONE (Supplementary Table 1). In other words, we still observed the general

386  pattern that nonsynonynmous variants were under-represented across the four

387     secondary structures and enriched over NONE, although at a smaller magnitude

388     compared to the overall analysis.

389

390     **Human intraspecific variation.** We identified 1,921 indels from 1,436 unique genes,

391     comprising a total of 4,354 positions. Most of these occurred at a frequency of 1 allele

392     observed among 5,008 phased alleles in the 1000 genomes project. We did not exclude

393     these; even if they are due to sequencing or mapping errors, there is no reason to

394     believe they would inflate our overall false positive rate as such errors should occur

395     blindly with respect to secondary structure of proteins. In addition, an indel at a

396     frequency of 1 allele could be especially deleterious, so we included them.

397        Across all 6 geographic regions, indel sites spanning NONE occurred at nearly

398     twice the frequency than secondary structures. NONE indels reached a mean frequency

399     of 4 alleles out of 5,008 phased alleles, compared to BEND/HELIX/TURN indels (3

400     alleles) and STRAND indels (1 allele) (Kruskal-Wallis $X^2$= 37.8, df = 2, p-value < $10^{-8}$). If

401     we use a minor allele frequency cutoff of 1%, 3% or 5% these patterns disappear,

402     indicating that the majority of signal comes from the fact that a large proportion of

403     STRAND indels occur as singletons.

404

405                         **Discussion**

406     Our study combined the recent revolution in protein structure, ushered in by the

407     AlphaFold2 project (Jumper et al. 2021), with evolutionary, population genetic, and

408     permutation-based analyses to demonstrate that indels were depleted in regions of

409     predicted secondary structure. This skew is especially strong for STRAND, which is

410    consistent with these structures being internal and stable regions that are important for

411    the overall 3D structure of a protein (Echave et al. 2016).

412         There are two non-mutually exclusive models – a mutational bias model versus a

413    selection model – that could explain the non-random distribution of indels that we

414    observe here. Under a mutational bias model, the four secondary structures experience

415    fundamentally different rates of indel mutation. The four different secondary structures

416    tested here display systematic differences in amino acid composition (Chou and

417    Fasman 1975; Fujiwara et al. 2012), which predicts different base composition and/or

418    repetitive elements in the underlying DNA, which in turn could influence mutation rate.

419         However, three patterns in our data argue against the mutational bias

420    hypothesis, and instead provide support for a model where selection acts against indels

421    that are more likely to disrupt protein function. First, *within* each secondary structure,

422    positions with indels tend to occur in externally oriented and high-pIDDT regions of

423    proteins (Fig. 4). A mutational bias hypothesis cannot account for this discrepancy

424    because they are the same secondary structures in different parts of the same protein.

425    Second, the observed vs. expected ratios (Table 1) are stronger in rodents compared to

426    primates (Figure 3). A mutational bias hypothesis cannot account for this interspecific

427    variation unless different species also experience different mutational biases. In

428    contrast, this pattern is predicted by a model of selection, because natural selection will

429    operate more efficiently in species with large effective population size (Kimura 1983;

430    Lynch 2007; Charlesworth 2009). Rodents have an effective population size that is

431    roughly 10-fold larger than primates (Ohta 1972; Zhao et al. 2000; Won and Hey 2005;

432    Geraldes et al. 2008; Geraldes et al. 2011). Finally, we showed that nonsynonymous

433    variants were also depleted in regions of secondary structure, although not to the same

434    degree (Table 1). A mutational bias hypothesis cannot explain the depletion of both

435    indels and nonsynonymous variants over secondary structure, because these two

436    classes differ in their mutational process.

437        To be sure, it is unlikely that indel mutations arise randomly. For example, G+C

438    content often correlates with a genomic region's susceptibility to insertions or deletions

439    (Sinden et al. 2002; Taylor et al. 2004), as well as features suggestive of a slippage

440    mechanism (Nishizawa and Nishizawa 2002). However, a model of selection does not

441    require indel mutation to be completely random. A selection model only requires any

442    non-randomness in mutational process to be equally distributed across the five

443    categories of secondary structure tested here. It should also be pointed out that our

444    study reports average deviations in observed vs. expected across the entire genome. It

445    remains unknown how much the strength of selection varies across individual indels,

446    although our Gene Ontology results did not uncover any functional similarity among the

447    most highly skewed genes.

448        It is noteworthy that even within humans, we observed proportionately fewest

449    indels over STRAND – exactly the secondary structure where indels were depleted in

450    our five species analyses. The low historical effective population size of humans,

451    coupled with multiple bottlenecks, are expected to reduce the efficiency of selection, yet

452    we still observe skews in indel locations.

453        In conclusion, our analyses indicate that any change in amino acid sequence is

454    likely to be deleterious for secondary structure, especially if that change is not a single

455    nonsynonymous mutation, but the insertion or deletion of multiple amino acids. Indels

456  that overlap STRAND and/or buried regions of the protein, appear to be the most

457  deleterious, while indels over NONE the least. By analyzing the AlphaFold2 predictions,

458  we have quantified these effects over whole genomes and full-length proteins, revealing

459  a role for protein structure on the evolution of its primary sequence.

460

461  **Data and resource availability**

462  All data, code, and intermediate files required to reproduce the results here, as well as a

463  README file, are available on Dryad (https://doi.org/10.5061/dryad.bk3j9kdk9) as a

464  single protein_structure.tar.gz file (8.5 Gb). [for reviewers only: that link is not yet public;

465  this link provides access:

466  https://datadryad.org/stash/share/5NLwY6IUt75oIgY16DgFRkywjBUx2eoela6RYDHFHd

467  g]

468  **Acknowledgements**

**Figure Legends**

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

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

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

500    difference in O:E in insertions and deletions, but the branches leading to rodent species

501    generally show stronger skews than branches leading to primates.

502

503    **Figure 4.** Weighted means of relative solvent accessibility (red, left axis) and plDDT

504    scores (blue, right axis) across secondary structures, stratified by sites occurring over

505    indels (X) versus sites not overlapping indels (O). Numbers on x axis indicate the

506    number of sites that overlap an indel versus not (separated by |).

507

508    **Figure 5.** Violin plot of the minor allele frequency of indels in protein coding regions,

509    segregating within humans, stratified by secondary structure. B/H/T = pooled

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

511    Figure includes all human populations pooled; results remain qualitatively the same if

512    we analyze populations separately.

513

514    **Supplementary Figure 1.** A repeat of Figure 2, but for each of the four different

515    subsets of indels.

516

517    **Supplementary Figure 3.** A repeat of Figure 3, but for each of the four different

518    subsets of indels.

519

520    **Supplementary Figure 3.** A repeat of Figure 4, but for each of the four different

521    subsets of indels.

522

523

# References

524     **References**

526    Arpino JAJ, Reddington SC, Halliwell LM, Rizkallah PJ, Jones DD. 2014. Random Single Amino
527         Acid Deletion Sampling Unveils Structural Tolerance and the Benefits of Helical Registry
528         Shift on GFP Folding and Structure. *Structure* 22:889–898.

529    Banerjee A, Levy Y, Mitra P. 2019. Analyzing Change in Protein Stability Associated with Single
530         Point Deletions in a Newly Defined Protein Structure Database. *J. Proteome Res.*
531         18:1402–1410.

532    Barton HJ, Zeng K. 2019. The Impact of Natural Selection on Short Insertion and Deletion
533         Variation in the Great Tit Genome. *Genome Biology and Evolution* 11:1514–1524.

534    Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful
535         approach to multiple testing. *Journal of the royal statistical society. Series B*
536         *(Methodological)* 57:289–300.

537    Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE.
538         2000. The Protein Data Bank. *Nucleic Acids Research* 28:235–242.

539    Bermejo-Das-Neves C, Nguyen H-N, Poch O, Thompson JD. 2014. A comprehensive study of
540         small non-frameshift insertions/deletions in proteins and prediction of their phenotypic
541         effects by a machine learning method (KD4i). *BMC Bioinformatics* 15:111.

542    Bustamante CD, Townsend JP, Hartl DL. 2000. Solvent Accessibility and Purifying Selection
543         Within Proteins of Escherichia coli and Salmonella enterica. *Molecular Biology and*
544         *Evolution* 17:301–308.

545    Byrska-Bishop M, Evani US, Zhao X, Basile AO, Abel HJ, Regier AA, Corvelo A, Clarke WE,
546         Musunuri R, Nagulapalli K, et al. 2022. High-coverage whole-genome sequencing of the
547         expanded 1000 Genomes Project cohort including 602 trios. *Cell* 185:3426-3440.e19.

548    Charlesworth B. 2009. Effective population size and patterns of molecular evolution and
549         variation. *Nat Rev Genet* 10:195–205.

550    de la Chaux N, Messer PW, Arndt PF. 2007. DNA indels in coding regions reveal selective
551         constraints on protein evolution in the human lineage. *BMC Evol Biol* 7:191.

552    Chen J, Guo J. 2021. Structural and functional analysis of somatic coding and UTR indels in
553         breast and lung cancer genomes. *Sci Rep* 11:21178.

554    Chong Z, Zhai W, Li C, Gao M, Gong Q, Ruan J, Li J, Jiang L, Lv X, Hungate E, et al. 2013. The
555         Evolution of Small Insertions and Deletions in the Coding Genes of Drosophila
556         melanogaster. *Molecular Biology and Evolution* 30:2699–2708.

557    Chou PY, Fasman GD. 1975. Conformational parameters for amino acids in helical, β-sheet, and
558        random coil regions calculated from proteins. *ACS Publications* [Internet]. Available
559        from: https://pubs.acs.org/doi/pdf/10.1021/bi00699a001

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

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

564    Echave J, Spielman SJ, Wilke CO. 2016. Causes of evolutionary rate variation among protein
565        sites. *Nat Rev Genet* 17:109–121.

566    Fitch WM, Smith TF. 1983. Optimal sequence alignments. *Proceedings of the National Academy
567        of Sciences* 80:1382–1386.

568    Franzosa EA, Xia Y. 2009. Structural Determinants of Protein Evolution Are Context-Sensitive at
569        the Residue Level. *Molecular Biology and Evolution* 26:2387–2395.

570    Fujiwara K, Toda H, Ikeguchi M. 2012. Dependence of α-helical and β-sheet amino acid
571        propensities on the overall protein fold type. *BMC Structural Biology* 12:18.

572    Gavrilov Y, Dagan S, Levy Y. 2015. Shortening a loop can increase protein native state entropy.
573        *Proteins: Structure, Function, and Bioinformatics* 83:2137–2146.

574    Gavrilov Y, Dagan S, Reich Z, Scherf T, Levy Y. 2018. An NMR Confirmation for Increased Folded
575        State Entropy Following Loop Truncation. *J. Phys. Chem. B* 122:10855–10860.

576    Geraldes A, Basset P, Gibson B, Smith KL, Harr B, Yu HT, Bulatova N, Ziv Y, Nachman MW. 2008.
577        Inferring the history of speciation in house mice from autosomal, X-linked, Y-linked and
578        mitochondrial genes. *Mol Ecol* 17:5349–5363.

579    Geraldes A, Basset P, Smith KL, Nachman MW. 2011. Higher differentiation among subspecies
580        of the house mouse (*Mus musculus*) in genomic regions with low recombination. *Mol
581        Ecol* 20:4722–4736.

582    Goldman N, Thorne JL, Jones DT. 1998. Assessing the Impact of Secondary Structure and
583        Solvent Accessibility on Protein Evolution. *Genetics* 149:445–458.

584    Gonzalez CE, Roberts P, Ostermeier M. 2019. Fitness Effects of Single Amino Acid Insertions and
585        Deletions in TEM-1 β-Lactamase. *Journal of Molecular Biology* 431:2320–2330.

586    Grantham R. 1974. Amino acid difference formula to help explain protein evolution. *Science*
587        185:862–864.

588    Grocholski T, Dinis P, Niiranen L, Niemi J, Metsä-Ketelä M. 2015. Divergent evolution of an
589        atypical S-adenosyl-l-methionine–dependent monooxygenase involved in anthracycline
590        biosynthesis. *Proceedings of the National Academy of Sciences* 112:9866–9871.

591    Guo B, Zou M, Wagner A. 2012. Pervasive Indels and Their Evolutionary Dynamics after the Fish-
592        Specific Genome Duplication. *Molecular Biology and Evolution* 29:3005–3022.

593    Halliwell LM, Jathoul AP, Bate JP, Worthy HL, Anderson JC, Jones DD, Murray JAH. 2018. ΔFlucs:
594        Brighter Photinus pyralis firefly luciferases identified by surveying consecutive single
595        amino acid deletion mutations in a thermostable variant. *Biotechnology and*
596        *Bioengineering* 115:50–59.

597    Hartl DL, Clark AG. 2007. Principles of population genetics. 4th ed. Sunderland, MA: Sinauer

598    Hedrick PW. 2005. Genetics of populations. 3rd ed. Boston: Jones and Bartlett

599    Hormozdiari F, Salari R, Hsing M, Schönhuth A, Chan SK, Sahinalp SC, Cherkasov A. 2009. The
600        Effect of Insertions and Deletions on Wirings in Protein-Protein Interaction Networks: A
601        Large-Scale Study. *Journal of Computational Biology* 16:159–167.

602    Hsing M, Cherkasov A. 2008. Indel PDB: A database of structural insertions and deletions
603        derived from sequence alignments of closely related proteins. *BMC Bioinformatics*
604        9:293.

605    Iengar P. 2012. An analysis of substitution, deletion and insertion mutations in cancer genes.
606        *Nucleic Acids Research* 40:6401–6413.

607    Jackson EL, Spielman SJ, Wilke CO. 2017. Computational prediction of the tolerance to amino-
608        acid deletion in green-fluorescent protein. *PLOS ONE* 12:e0164905.

609    Jayaraman V, Toledo-Patiño S, Noda-García L, Laurino P. 2022. Mechanisms of protein
610        evolution. *Protein Science* 31:e4362.

611    Jilani M, Haspel N, Jagodzinski F. 2022. Detection and Analysis of Amino Acid Insertions and
612        Deletions. In: Haspel N, Jagodzinski F, Molloy K, editors. Algorithms and Methods in
613        Structural Bioinformatics. Computational Biology. Cham: Springer International
614        Publishing. p. 89–99. Available from: https://doi.org/10.1007/978-3-031-05914-8_5

615    Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates R,
616        Žídek A, Potapenko A. 2021. Highly accurate protein structure prediction with
617        AlphaFold. *Nature* 596:583–589.

618    Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple
619        sequence alignment based on fast Fourier transform. *Nucleic acids research* 30:3059–
620        3066.

Khan T, Douglas GM, Patel P, Nguyen Ba AN, Moses AM. 2015. Polymorphism Analysis Reveals Reduced Negative Selection and Elevated Rate of Insertions and Deletions in Intrinsically Disordered Protein Regions. *Genome Biology and Evolution* 7:1815–1826.

Kim PM, Lu LJ, Xia Y, Gerstein MB. 2006. Relating Three-Dimensional Structures to Protein Networks Provides Evolutionary Insights. *Science* 314:1938–1941.

Kim R, Guo J. 2010. Systematic analysis of short internal indels and their impact on protein folding. *BMC Struct Biol* 10:24.

Kimura M. 1983. The neutral theory of molecular evolution. Cambridge: Cambridge University Press

Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH, Ware JS, Hill AJ, Cummings BB, et al. 2016. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 536:285–291.

Levy Karin E, Susko E, Pupko T. 2014. Alignment Errors Strongly Impact Likelihood-Based Tests for Comparing Topologies. *Molecular Biology and Evolution* 31:3057–3067.

Light S, Sagit R, Ekman D, Elofsson A. 2013. Long indels are disordered: A study of disorder and indels in homologous eukaryotic proteins. *Biochimica et Biophysica Acta (BBA) - Proteins and Proteomics* 1834:890–897.

Light S, Sagit R, Sachenkova O, Ekman D, Elofsson A. 2013. Protein Expansion Is Primarily due to Indels in Intrinsically Disordered Regions. *Molecular Biology and Evolution* 30:2645–2653.

Lin M, Whitmire S, Chen J, Farrel A, Shi X, Guo J. 2017. Effects of short indels on protein structure and function in human genomes. *Sci Rep* 7:9313.

Liu S, Wei X, Dong X, Xu L, Liu J, Jiang B. 2015. Structural plasticity of green fluorescent protein to amino acid deletions and fluorescence rescue by folding-enhancing mutations. *BMC Biochemistry* 16:17.

Liu S, Wei X, Ji Q, Xin X, Jiang B, Liu J. 2016. A facile and efficient transposon mutagenesis method for generation of multi-codon deletions in protein sequences. *Journal of Biotechnology* 227:27–34.

Lynch M. 2007. The origins of genome architecture. Available from: https://repository.library.georgetown.edu/handle/10822/548280

Mariani V, Biasini M, Barbato A, Schwede T. 2013. lDDT: a local superposition-free score for comparing protein structures and models using distance difference tests. *Bioinformatics* 29:2722–2728.

654    Marsh JA, Teichmann SA. 2014. Parallel dynamics and evolution: Protein conformational
655         fluctuations and assembly reflect evolutionary changes in sequence and structure.
656         *BioEssays* 36:209–218.

657    Mi H, Huang X, Muruganujan A, Tang H, Mills C, Kang D, Thomas PD. 2017. PANTHER version 11:
658         expanded annotation data from Gene Ontology and Reactome pathways, and data
659         analysis tool enhancements. *Nucleic acids research* 45:D183–D189.

660    Mi H, Muruganujan A, Ebert D, Huang X, Thomas PD. 2019. PANTHER version 14: more
661         genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools.
662         *Nucleic acids research* 47:D419–D426.

663    Mi H, Muruganujan A, Thomas PD. 2013. PANTHER in 2013: modeling the evolution of gene
664         function, and other gene attributes, in the context of phylogenetic trees. *Nucleic Acids*
665         *Res* 41:D377-386.

666    Mills RE, Pittard WS, Mullaney JM, Farooq U, Creasy TH, Mahurkar AA, Kemeza DM, Strassler
667         DS, Ponting CP, Webber C, et al. 2011. Natural genetic variation caused by small
668         insertions and deletions in the human genome. *Genome Res.* 21:830–839.

669    Mintseris J, Weng Z. 2005. Structure, function, and evolution of transient and obligate protein–
670         protein interactions. *Proceedings of the National Academy of Sciences* 102:10930–
671         10935.

672    Mitternacht S. 2016. FreeSASA: An open source C library for solvent accessible surface area
673         calculations. *F1000Res* 5:189.

674    Montgomery SB, Goode DL, Kvikstad E, Albers CA, Zhang ZD, Mu XJ, Ananda G, Howie B,
675         Karczewski KJ, Smith KS, et al. 2013. The origin, evolution, and functional impact of short
676         insertion–deletion variants identified in 179 human genomes. *Genome Res.* 23:749–761.

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

679    Nishizawa M, Nishizawa K. 2002. A DNA Sequence Evolution Analysis Generalized by Simulation
680         and the Markov Chain Monte Carlo Method Implicates Strand Slippage in a Majority of
681         Insertions and Deletions. *J Mol Evol* 55:706–717.

682    Ohta T. 1972. Evolutionary rate of cistrons and DNA divergence. *J Mol Evol* 1:150–157.

683    Pascarella S, Argos P. 1992. Analysis of insertions/deletions in protein structures. *Journal of*
684         *Molecular Biology* 224:461–471.

685    Penn O, Privman E, Ashkenazy H, Landan G, Graur D, Pupko T. 2010. GUIDANCE: a web server
686         for assessing alignment confidence scores. *Nucleic acids research* 38:W23–W28.

687     Penn O, Privman E, Landan G, Graur D, Pupko T. 2010. An alignment confidence score capturing
688          robustness to guide tree uncertainty. *Mol Biol Evol* 27:1759–1767.

689     Privman E, Penn O, Pupko T. 2012. Improving the performance of positive selection inference
690          by filtering unreliable alignment regions. *Molecular Biology and Evolution* 29:1–5.

691     Rockah-Shmuel L, Tóth-Petróczy Á, Sela A, Wurtzel O, Sorek R, Tawfik DS. 2013. Correlated
692          Occurrence and Bypass of Frame-Shifting Insertion-Deletions (InDels) to Give Functional
693          Proteins. *PLOS Genetics* 9:e1003882.

694     Salari R, Schönhuth A, Hormozdiari F, Cherkasov A, Sahinalp SC. 2008. The Relation between
695          Indel Length and Functional Divergence: A Formal Study. In: Crandall KA, Lagergren J,
696          editors. Algorithms in Bioinformatics. Lecture Notes in Computer Science. Berlin,
697          Heidelberg: Springer. p. 330–341.

698     Savino S, Desmet T, Franceus J. 2022. Insertions and deletions in protein evolution and
699          engineering. *Biotechnology Advances* 60:108010.

700     Scherrer MP, Meyer AG, Wilke CO. 2012. Modeling coding-sequence evolution within the
701          context of residue solvent accessibility. *BMC Evolutionary Biology* 12:179.

702     Shahmoradi A, Sydykova DK, Spielman SJ, Jackson EL, Dawson ET, Meyer AG, Wilke CO. 2014.
703          Predicting Evolutionary Site Variability from Structure in Viral Proteins: Buriedness,
704          Packing, Flexibility, and Design. *J Mol Evol* 79:130–142.

705     Shih C-H, Chang C-M, Lin Y-S, Lo W-C, Hwang J-K. 2012. Evolutionary information hidden in a
706          single protein structure. *Proteins: Structure, Function, and Bioinformatics* 80:1647–1657.

707     Simm AM, Baldwin AJ, Busse K, Jones DD. 2007. Investigating protein structural plasticity by
708          surveying the consequence of an amino acid deletion from TEM-1 β-lactamase. *FEBS
709          Letters* 581:3904–3908.

710     Sinden RR, Potaman VN, Oussatcheva EA, Pearson CE, Lyubchenko YL, Shlyakhtenko LS. 2002.
711          Triplet repeat DNA structures and human genetic disease: dynamic mutations from
712          dynamic DNA. *J Biosci* 27:53–65.

713     Slodkowicz G, Goldman N. 2020. Integrated structural and evolutionary analysis reveals
714          common mechanisms underlying adaptive evolution in mammals. *Proc Natl Acad Sci
715          USA* 117:5977–5986.

716     Snir S, Pachter L. 2006. Phylogenetic Profiling of Insertions and Deletions in Vertebrate
717          Genomes. In: Apostolico A, Guerra C, Istrail S, Pevzner PA, Waterman M, editors.
718          Research in Computational Molecular Biology. Lecture Notes in Computer Science.
719          Berlin, Heidelberg: Springer. p. 265–280.

720 Tao S, Fan Y, Wang W, Ma G, Liang L, Shi Q. 2007. Patterns of Insertion and Deletion in
721     Mammalian Genomes. *Current Genomics* 8:370–378.

722 Taylor MS, Ponting CP, Copley RR. 2004. Occurrence and Consequences of Coding Sequence
723     Insertions and Deletions in Mammalian Genomes. *Genome Res.* 14:555–566.

724 The Genomes Project C. 2015. A global reference for human genetic variation. *Nature* 526:68–
725     74.

726 Thomas PD, Ebert D, Muruganujan A, Mushayahama T, Albou L-P, Mi H. 2022. PANTHER:
727     Making genome-scale phylogenetics accessible to all. *Protein Science* 31:8–22.

728 Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. 2013. Maximum Allowed Solvent
729     Accessibilites of Residues in Proteins. *PLoS One* 8:e80635.

730 Tóth-Petróczy Á, Tawfik DS. 2011. Slow protein evolutionary rates are dictated by surface–core
731     association. *Proceedings of the National Academy of Sciences* 108:11151–11156.

732 Tóth-Petróczy Á, Tawfik DS. 2014. Hopeful (Protein InDel) Monsters? *Structure* 22:803–804.

733 Varadi M, Anyango S, Deshpande M, Nair S, Natassia C, Yordanova G, Yuan D, Stroe O, Wood G,
734     Laydon A. 2022. AlphaFold Protein Structure Database: massively expanding the
735     structural coverage of protein-sequence space with high-accuracy models. *Nucleic acids*
736     *research* 50:D439–D444.

737 Won YJ, Hey J. 2005. Divergence population genetics of chimpanzees. *Mol Biol Evol* 22:297–307.

738 Woods H, Schiano DL, Aguirre JI, Ledwitch KV, McDonald EF, Voehler M, Meiler J, Schoeder CT.
739     2023. Computational modeling and prediction of deletion mutants. *Structure* 31:713-
740     723.e3.

741 Yeh S-W, Liu J-W, Yu S-H, Shih C-H, Hwang J-K, Echave J. 2014. Site-Specific Structural
742     Constraints on Protein Sequence Evolutionary Divergence: Local Packing Density versus
743     Solvent Exposure. *Molecular Biology and Evolution* 31:135–139.

744 Zhang Z, Huang J, Wang Z, Wang L, Gao P. 2011. Impact of Indels on the Flanking Regions in
745     Structural Domains. *Molecular Biology and Evolution* 28:291–301.

746 Zhang Z, Wang J, Gong Y, Li Y. 2018. Contributions of substitutions and indels to the structural
747     variations in ancient protein superfamilies. *BMC Genomics* 19:771.

748 Zhang Z, Wang Y, Wang L, Gao P. 2010. The Combined Effects of Amino Acid Substitutions and
749     Indels on the Evolution of Structure within Protein Families. *PLOS ONE* 5:e14316.

750     Zhao Z, Jin L, Fu YX, Ramsay M, Jenkins T, Leskinen E, Pamilo P, Trexler M, Patthy L, Jorde LB, et

751         al. 2000. Worldwide DNA sequence variation in a 10-kilobase noncoding region on

752         human chromosome 22. *Proc Natl Acad Sci USA* 97:11354–11358.
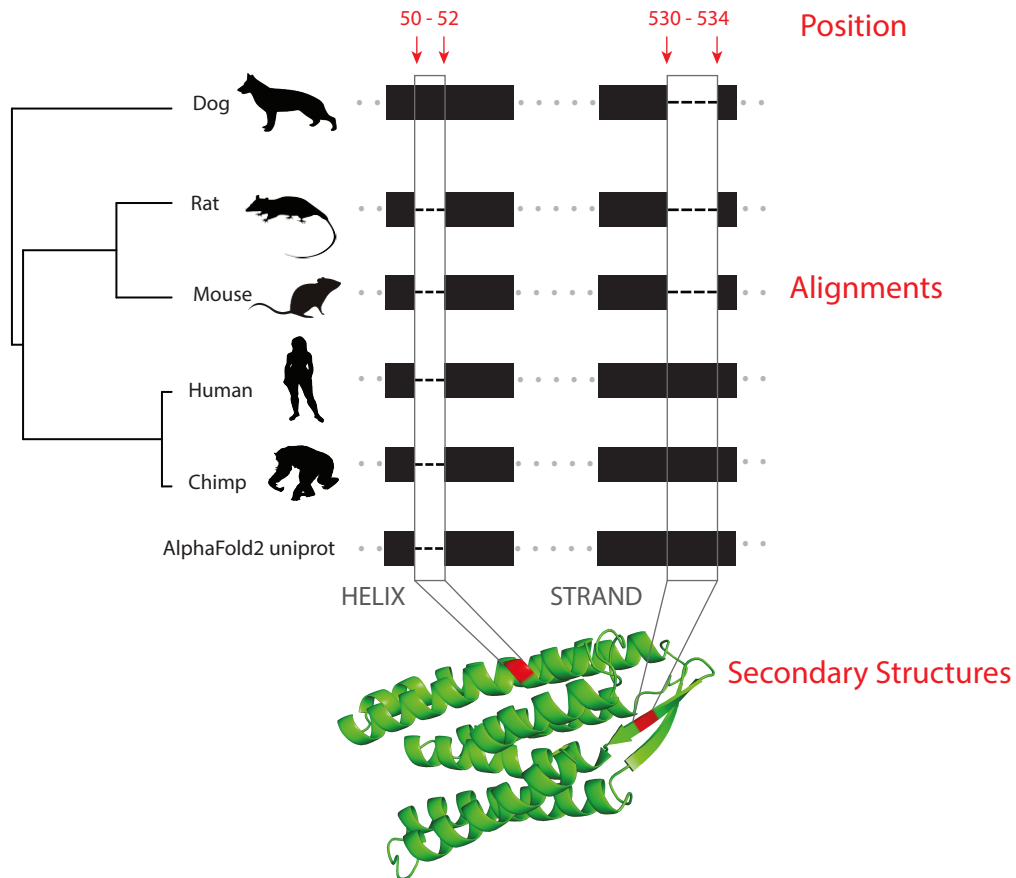
753

754

755

Figure 1

Figure 2



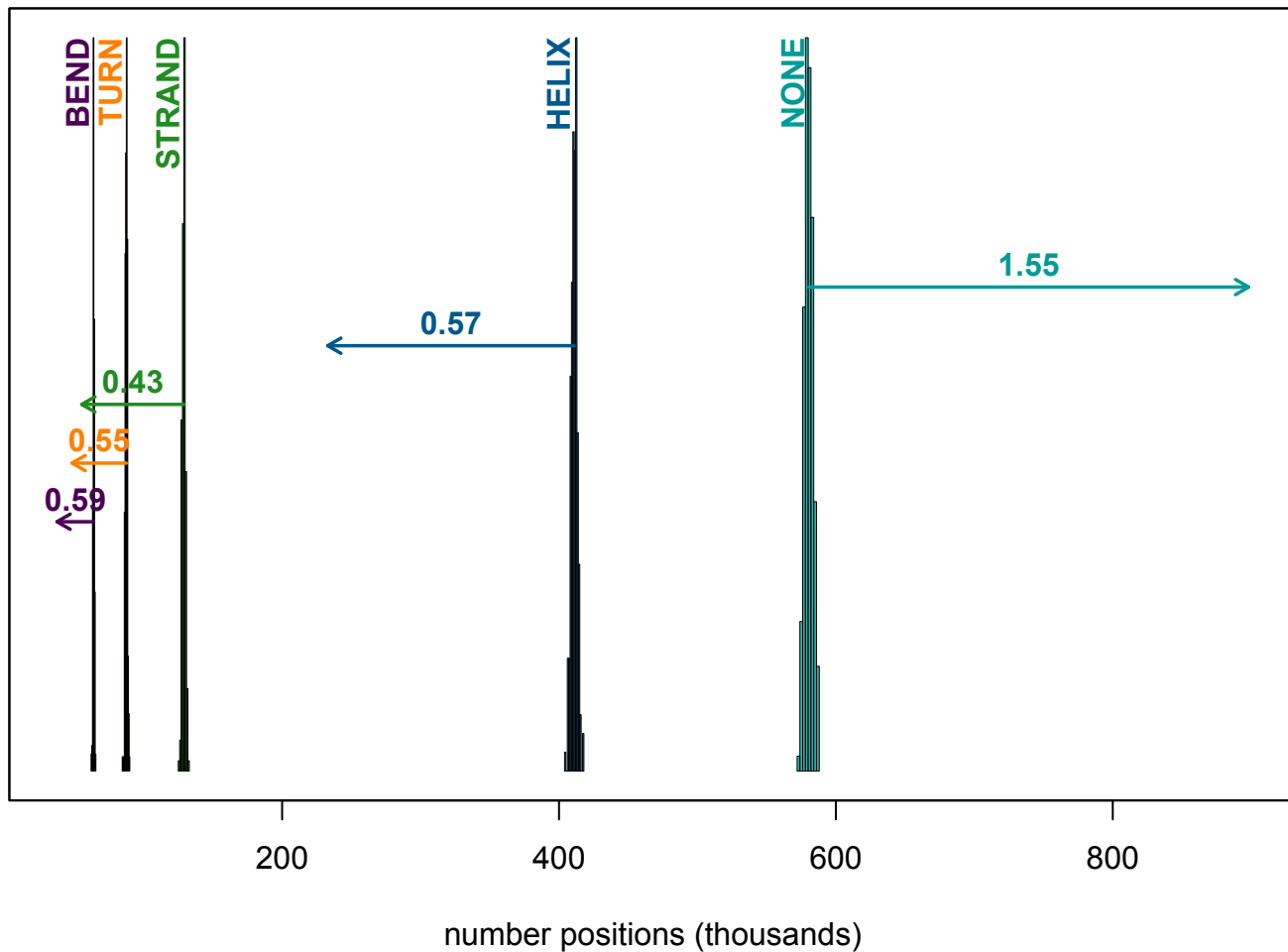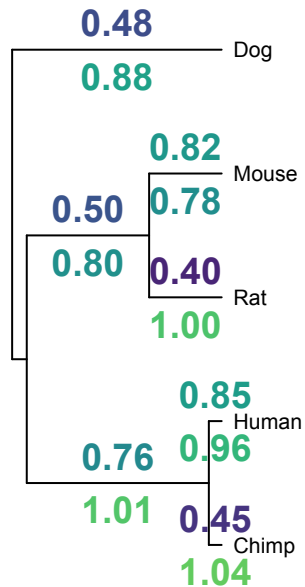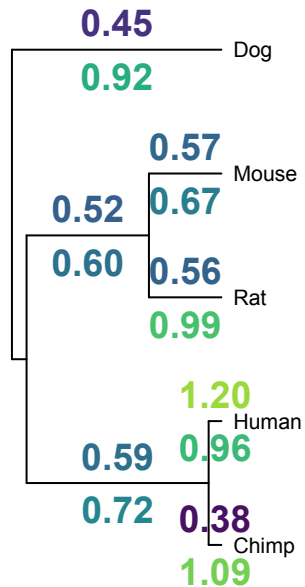**ALL**
**11444 genes, 97383 indels, 200 iterations**

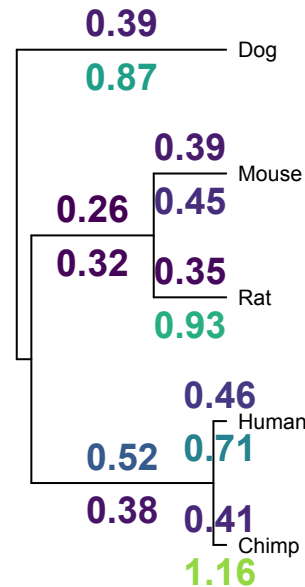number positions (thousands)

Figure 3

**BEND**

- 0.48
- Dog
- 0.88
- 0.82
- 0.50
- 0.78
- Mouse
- 0.80
- 0.40
- Rat
- 1.00
- 0.85
- 0.76
- 0.96
- Human
- 1.01
- 0.45
- Chimp
- 1.04

**TURN**

- 0.45
- Dog
- 0.92
- 0.57
- 0.52
- 0.67
- Mouse
- 0.60
- 0.56
- Rat
- 0.99
- 1.20
- 0.59
- 0.96
- Human
- 0.72
- 0.38
- Chimp
- 1.09

**STRAND**

- 0.39
- Dog
- 0.87
- 0.39
- 0.26
- 0.45
- Mouse
- 0.32
- 0.35
- Rat
- 0.93
- 0.46
- 0.52
- 0.71
- Human
- 0.38
- 0.41
- Chimp
- 1.16

**HELIX**

- 0.48
- Dog
- 0.91
- 0.67
- 0.46
- 0.65
- Mouse
- 0.52
- 0.49
- Rat
- 0.96
- 0.82
- 0.58
- 0.95
- Human
- 0.60
- 0.47
- Chimp
- 1.04

**NONE**

- 1.64
- Dog
- 1.15
- 1.52
- 1.65
- 1.41
- Mouse
- 1.50
- 1.73
- Rat
- 1.06
- 1.30
- 1.57
- 1.12
- Human
- 1.45
- 1.79
- Chimp
- 0.89
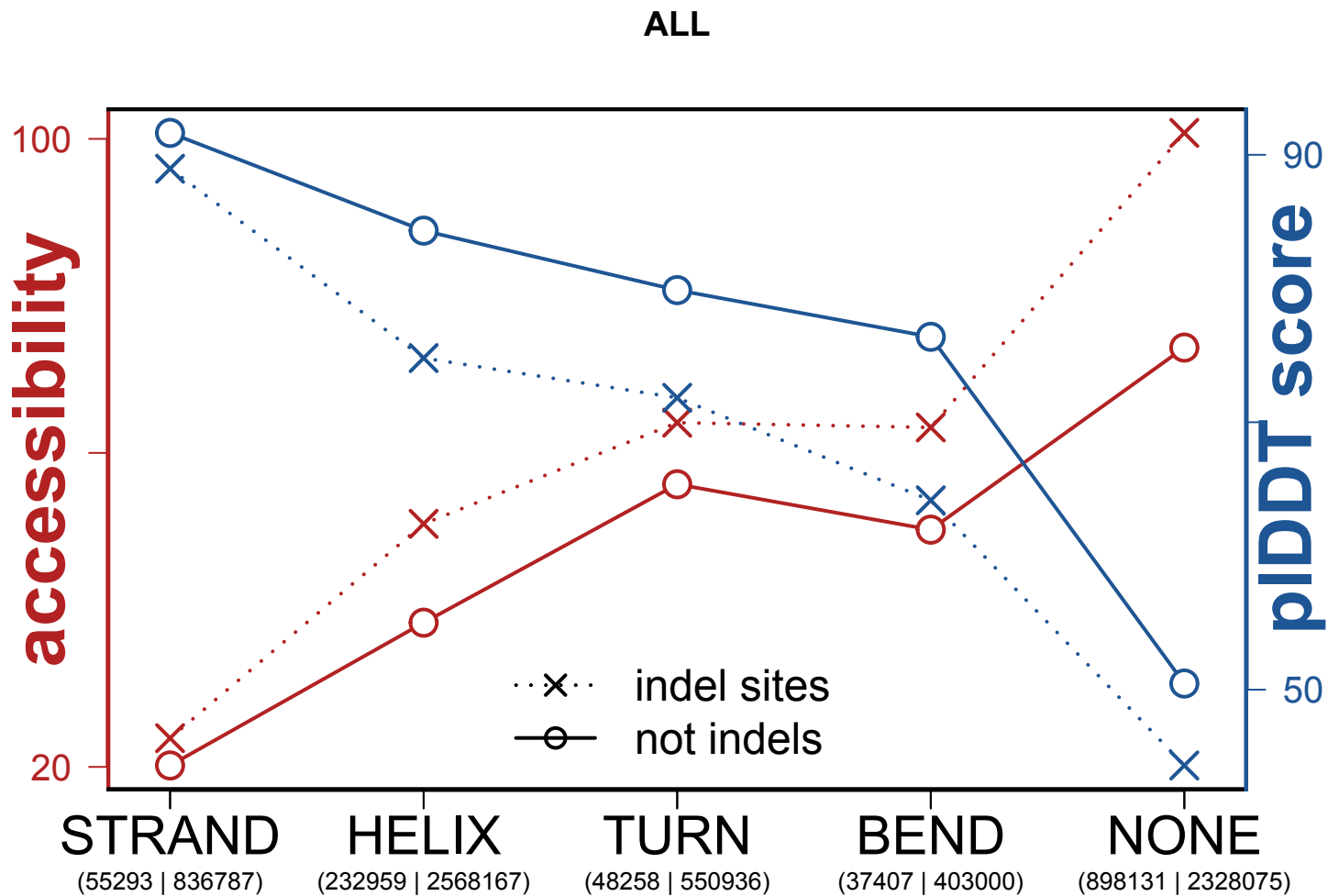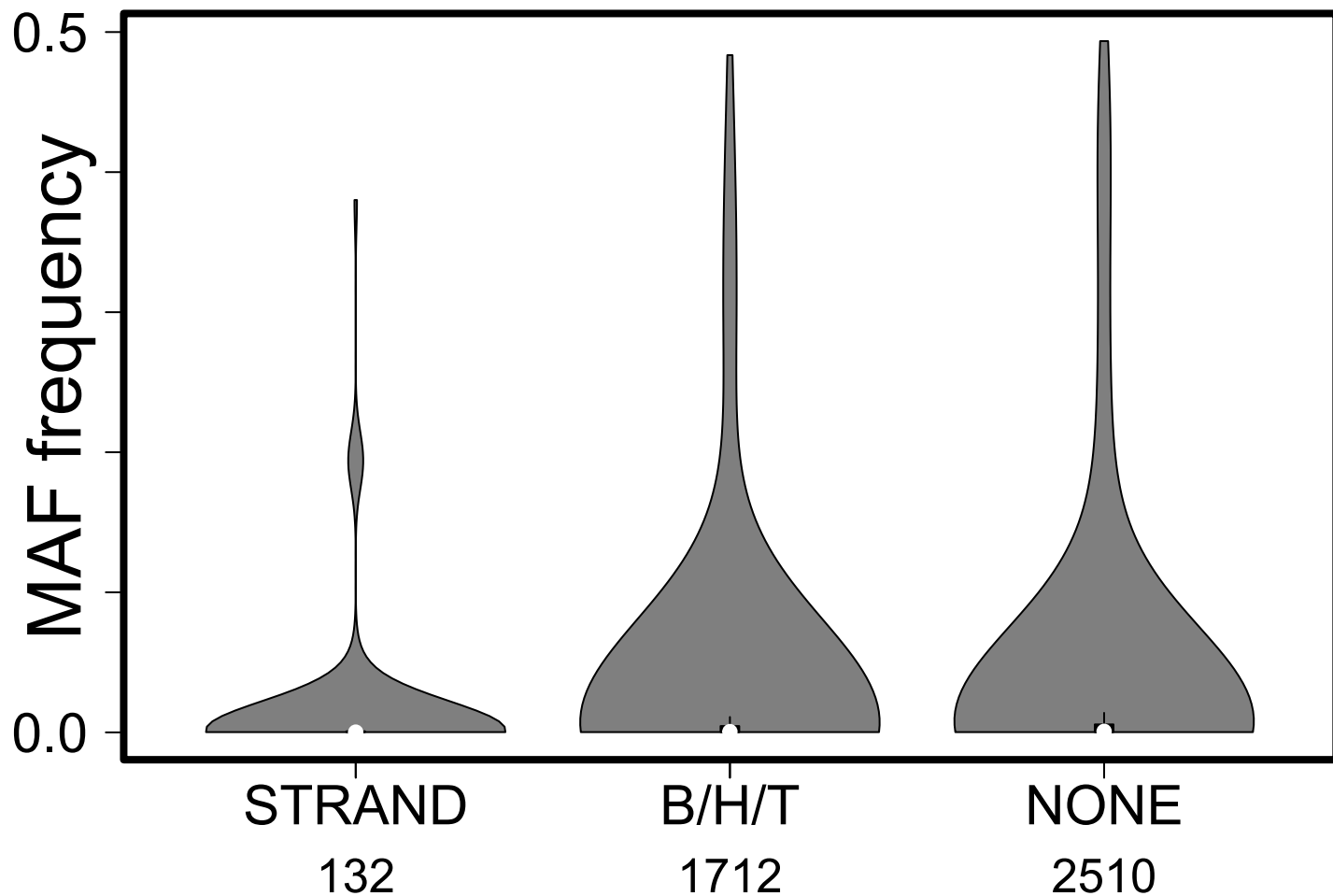
ALL

Figure 4



ALL

Figure 5

**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 nonsynonsymous (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) |