# CAUSA 2.0: accurate and consistent evolutionary analysis of proteins using codon and amino acid unified sequence alignments

Xiaolong Wang, Chao Yang

Multiple sequence alignment (MSA) is widely used to reveal structural and functional changes leading to genetic differences among species, and to reconstruct evolutionary histories of related genes, proteins and genomes. Traditionally, proteins and their coding sequences (CDSs) are aligned and analyzed separately, but often drastically different conclusions were drawn on a same set of data. Here we present a new alignment strategy, *Codon and Amino Acid Unified Sequence Alignment (CAUSA) 2.0*, which aligns proteins and their coding sequences simultaneously. CAUSA 2.0 optimizes the alignment of CDSs at both codon and amino acid level efficiently. Theoretical analysis showed that CAUSA 2.0 enhances the entropy information content of MSA. Empirical data analysis demonstrated that CAUSA 2.0 is more accurate and consistent than nucleotide, protein or codon level alignments. CAUSA 2.0 locates in-frame indels more accurately, makes the alignment of coding sequences biologically more significant, and reveals several novel mutation mechanisms that relate to some genetic diseases. CAUSA 2.0 is available in website www.DNAPlusPro.com .

**CAUSA 2.0: accurate and consistent evolutionary analysis of proteins using codon and amino acid unified sequence alignments**

*Xiaolong Wang*, Chao Yang*

*Department of Biotechnology, Ocean University of China, Qingdao, 266003, China*

*Correspondence:

*Xiaolong Wang, Ph.D., Associate Professor,*

*Department of Biotechnology*

*Ocean University of China*

*No. 5 Yushan Road*

*Qingdao, 266003, Shandong, P. R. China*

*Tel: 0086-139-6969-3150*

*E-mail: Xiaolong@ouc.edu.cn*

## Introduction

In recent years, with increasingly wider availability of sequence data, it has been able to reveal structural and functional changes leading to genetic differences among species, and allows reconstruction of evolutionary histories of related genes, proteins and genomes. These studies require multiple sequence alignment (MSA) for accurate identification of homologous residues and positioning of gaps indicating insertions and deletions (indels). At present, there are quite some MSA tools, such as Clustal W [1,2], MAFFT [3], MUSCLE [4], T-coffee [5] and PRANK [6]. However, different programs often give drastically different alignments on a same set of data, and support entirely different mechanisms driven evolutionary changes. In fact, MSAs are sometimes so inaccurate that they could lead to misinterpretation of evolutionary information [7].

Due to the small size of the alphabet of DNA bases, it is difficult to distinguish true homology from random similarities, and thus the alignment of DNA sequences is inherently inaccurate [8]. Proteins are built from 20 amino acids, so that the 'signal-to-noise ratio' in protein sequence alignments is better than that of DNA sequences. Besides this advantage in information-content, protein alignments also benefit from a series of AA substitution matrices, such as PAM [9], BLOSUM [10] and Gonnet [11], which contain empirically derived scores for each possible AA substitution and provide a rational basis for aligning AAs.

Particularly, the problem of aligning the coding sequences (CDSs) for a certain family of proteins has been studied fairly in depth. Various methods have been developed to solve this problem: The first is to construct a codon alignment by back-translating a protein alignment, such as RevTrans [8], PAL2NAL [12] and TranslatorX [13]; The second, Hein's COMBAT [14,15], is to construct a combined alignment by combining a DNA alignment and a protein alignment; Unfortunately, it has been applied only in pairwise alignment, but not in MSA; The third is to construct a codon alignment by using an empirical codon substitution scoring matrix [16]. In addition, many codon-based models [17-21] has been developed to produce more reliable estimates of certain biologically important measures, such as the transition/transversion rate ratio and the synonymous/nonsynonymous substitution rate ratio, than can nucleotide-based models.

43    Owing to the complexities of the alignment of CDSs, it has been rather difficult to choose an

44    appropriate tool, method and parameters for MSA and the subsequent evolutionary analysis.

45    Moreover, traditionally protein sequences and their CDSs have been aligned and analyzed

46    separately, in nucleotide, amino acid (AA) or codon level. It is well known that they often give

47    incongruent conclusions, in practice, however, usually only one kind of alignment is selected for

48    the subsequent phylogenetic and evolutionary analysis, while the others are dismissed to avoid

49    the problem of inconsistency in different (nucleotide, AA or codon) alignment methods. Here we

50    present an alternative alignment strategy, *Codon and Amino Acid Unified Sequence Alignment*

51    (CAUSA 2.0), which aligns protein sequences and their CDSs simultaneously. We demonstrated

52    that CAUSA 2.0 is evolutionarily accurate and consistent, computationally efficient, and it

53    reveals some novel mutation mechanisms.


54    **Materials and Methods**


55    *Protein coding sequences*

56    Different strains of human and simian immunodeficiency virus were derived from the seed

57    alignment of Pfam family pf00516. The coding sequences (CDSs) of their Envelope glycoprotein

58    gp120 (Env) and core (Gag) proteins were retrieved from the HIV sequence database [22]. Thirty

59    protein families and their standard phylogenetic trees in human and mammalian animals were

60    arbitrarily selected from TreeFam-A (http://www.treefam.org/).

61    *Data analysis*

62    The flowchart of data analysis is shown in Fig 1, DNA or protein alignments were aligned

63    by the multiple sequence alignment tools at EBI (http://www.ebi.ac.uk/), including ClustalW,

64    MAFFT, MUSCLE, T-COFFEE and PRANK. Codon alignments were aligned by PRANK using

65    "align translated codons" option (PRANK-CA), and a codon alignment tool (CAT) provided by

66    the HIV database at the Los Alamos National Laboratory (http://www.hiv.lanl.gov/). All

67    programs were run with their default settings. The principle and implementation of CAUSA 2.0,

68    and the methodologies for the construction of sequence alignments, phylogenetic trees,

69   evolutionary analysis and structure modeling, were described in details in the *Supplementary*

70   *Material*.


## 71   Results


### 72   *CAUSA 2.0 improves the accuracy of the alignment of CDSs*

73   It has been reported that traditional nucleotide or protein level alignment methods

74   incorrectly squeeze distinct, but similar, inserted sequences between two conserved blocks, thus

75   overestimate the number of substitutions and underestimate that of indels [23]. Using the Pfam seed

76   alignment of Env (PF00516) as a model example, as shown in Fig 2A and Fig S1A, they

77   demonstrated that a traditional ClustalW protein alignment of Env suggested a high rate of

78   substitutions in a variable (V2) region. Alignments of the same region given by MAFFT (Fig

79   S1B), MUSCLE (Fig S1C) and T-coffee (Fig S1D), are basically the same to that of ClustalW.

80   They also pointed out that the problem is caused by repeated penalizing gap-opening [23], but

81   cannot be avoided by reducing gap-opening penalties, since it will result in 'gappy' alignments.

82   In order to solve this problem, they developed PRANK [6,7,23-26], a phylogeny-aware alignment

83   tool, which "flags" gaps introduced in earlier steps, so that distinct insertions are kept separate

84   even when they occur at a same position [23]. As shown in Fig 2B, PRANK identified several

85   distinct insertions.

86   Using the same set of Env CDSs, we compared CAUSA 2.0 with protein alignments aligned

87   by ClustalW, MAFFT, MUSCLE, T-COFFEE and PRANK, and codon alignments aligned by

88   PRANK-Codon and CAT. As shown in Fig S1A-S1D, all different kinds of protein alignments

89   show a lot of AA substitutions in the V2 region of Env, while in PRANK (Fig 2B, S1E) and

90   PRANK-Codon (Fig 2C, S1F), they are identified as distinct insertions, and thus, many more

91   gaps were inserted.

92   In CAUSA 2.0 non-synonymous substitutions (NSSs) are displayed in obviously different

93   colors; and synonymous substitutions (SSs) are shown in similar but slightly different colors,

94   makes it much easier for a user to distinguish an SS from a NSS. As indicated by the red boxes

95  in Fig 2, in all of the other alignments, a row of highly conserved Y residues were misaligned.

96  While in CAUSA 2.0 the alignment of this row is obviously more accurate. It is more clearly

97  shown in Fig S1H, the CAUSA 2.0 alignment is visually highly ordered, since the number of

98  substitutions was minimized both in codon and AA level. Similar results could not be obtained

99  with the other MSA by changing their gap penalties. Alignments created by the other programs

100  often require fine adjustments made by human visual inspection before phylogeny analysis,

101  which is cumbersome and a potential source of errors. Using CAUSA 2.0, the need of visual

102  control or manual editing of MSAs is eliminated.

103  ***CAUSA 2.0 better interprets molecular evolution***

104  Frequently, CAUSA 2.0 inserts one or more gaps within a codon-AA 4-tuple. In the

105  previous version (CAUSA 1.0), we took them as errors, and had to fix the alignments a posteriori

106  by moving the gaps inserted out of the tuple [27]. However, later we found that gap-moving is

107  totally unnecessary. In the present version (CAUSA 2.0), the gapped codon-AA tuples are kept

108  and called *split codons*. Actually, the split codons are very informative and useful, as they allow

109  CAUSA 2.0 locates in-frame indels more accurately, and reveals some novel mutation

110  mechanisms.

111  An in-frame indel is an indel in a CDS with a length of a multiple of 3bp, and thus,

112  maintained the open reading frame (ORF). It has been reported that more than half of the

113  functionally important coding indels in human genome were in-frame, but they often do not

114  coincide perfectly with codon boundaries [28]. However, in a traditional protein or codon

115  alignment, every codon is indivisible, and thus, an in-frame indel has to be forced to start from

116  the first base of a codon, leads to an inaccurate alignment of the CDSs involved. By allowing the

117  split codons, CAUSA 2.0 locates in-frame and non-in-frame indels more accurately, and reveals

118  several new mutation mechanisms by which new codons could be generated:

119  (1) *Codon splitting*: when an in-frame insertion was inserted into an internal location

120  (between 1-2 or 2-3) of a codon, the codon was split into two parts, and fused with the bases in

121  the two ends of the inserted sequence, respectively, forming two new codons.

122    (2) *Codon fusion*: when an in-frame deletion started from the second or third base of a

123    codon, the remaining one or two bases of this codon fused with the bases in the other side of the

124    deletion, forming a new codon.

125    (3) *Substitution induced deletion* (*siDel*): when a non-synonymous substitution occurred in

126    a codon, if the encoded AA became similar to another AA nearby in the same sequence, it may

127    allow the deletion of a fragment between them.

128    (4) *indel-induced partial frame shift* (*ipFS*): when two or more non-in-frame indels

129    neighbored, there will be one or more partial frame shifts occurred within or nearby the indels, if

130    it does not change the whole reading frame of the CDSs.

131    For an example of codon splitting, as shown in Fig 3A and 3B, ClustalW suggested a lot of

132    independent substitutions between the two groups of virus (HIV1 and HIV2/SIV), PRANK-

133    Codon shows only a few indels and substitutions. However, as shown in Fig 3C, CAUSA 2.0

134    split codon 670 (Ggg-g) and located the 6-nt insertion (a gca ct) between the 2nd and the 3rd

135    base of this codon. Comparing the three alignments, CAUSA 2.0 might have the highest

136    probability of correctness, since it has the highest number of matches, and requires the least

137    number of mutation events.

138    In all HIV strains, HIV-1 subtypes are the most common form, they are the most infectious

139    and pathogenic to human, and among the most genetically diverse. During viral penetration into

140    host cell, the fusion peptide region involves in the merging of the virus envelop with host

141    endosomal membrane. The inserted motif is highly conserved in all HIV-1 strains, and is the

142    only difference located within the fusion peptide region compared with the other strains.

143    Therefore, we postulate that the insertion might be a key event lead to the origin of HIV-1, and a

144    critical motif in the interaction of the virus with human cell during the infection process.

145    For an example of codon fusion, as shown in Fig 3D, ClustalW shows a lot of substitutions,

146    PRANK-Codon shows that there are several indels (Fig 3E). CAUSA 2.0 (Fig 3F), however,

147    identified some smaller, but more accurate, indels. In addition, in different strains the codons

148    between 864 and 867 were deleted and mutated in different ways, however, the amino acid

149 encoded by codon 864 (V or C), is highly conserved among all strains, so it might probably be
150 functionally important.

151     Note that although a codon spitting could be distinguished easily from a codon fusion
152 according to a phylogenetic tree, it is hard to distinguish an ipFS from a siDel if their ancestral
153 sequence is not available. For example, as shown in Fig 3G and 3H, in ClustalW there are so
154 many base and AA substitutions that PRANK-Codon suggested long indels instead. CAUSA 2.0,
155 however, revealed multiple mutational events, including codon fusion/splitting, siDel and/or
156 ipFS, and gave a different while more accurate alignment for this region (Fig 3I). These
157 evolutionary events bring additional variability, in together with those by the traditional mutation
158 mechanisms, help the HIV to evolve rapidly and develop the ability of drug resistance.

159     By analyzing the CDSs of thirty protein families in human and mammalian arbitrarily
160 selected from TreeFam-A (Table S3), we found that these new mutation events are not rare, but
161 happened rather frequently not only in virus, but also in human and animals. We postulate that
162 these newly discovered mutation mechanisms may play an important role in the acquisition of
163 new genetic information and provide a potentially powerful way to optimize the protein
164 sequences and structures.

165 ***CAUSA 2.0 better explains certain genetic diseases***

166     Over the past decade, millions of small indels have been discovered in human populations
167 and personal genomes [28]. Most of the indels found in protein coding sequences were codon-
168 sized, in-frame indels, and thus maintained the ORF of these proteins. Numerous indels map to
169 functionally important sites within human genes, and are related to human traits and genetic
170 diseases [29-34].

171     For example, in-frame indels have been identified as one of the main causes for the Wolfram
172 syndrome, an autosomal recessive disorder [35]. Table 1 lists some in-frame indels related to
173 Wolfram syndrome and other genetic diseases, which were deposit in the Human Gene Mutation
174 Database. For example (CD031549), in the coding sequence …K*a*^**ag Agca Ag**^*cc*…, after the
175 deletion of the six nucleotides (in bold), the remaining three nucleotides (in italic) fused into a

176 new codon, Tacc, in which an amino acid substitution occurred without any base substitution. As

177 shown in Fig 4, such kind of disease-related codon splitting/fusion phenomenon is visualized

178 clearly in CAUSA 2.0, but not in the other alignment software, although a nucleotide alignment

179 may sometimes result in the same alignment. Therefore, CAUSA 2.0 would be useful in

180 discovering, recognizing and investigation of this kind of genetic diseases

181 ***CAUSA 2.0 improves molecular phylogenetic analyses in virus***

182 For several reasons ranging from methodological issues or bona fide biological phenomena,

183 the phylogenetic trees of different genes for a set of certain genomes are often incongruent

184 topologically. In order to tackle the topological conflicts in different gene trees, phylogenomic

185 studies couple concatenation with practices such as rogue taxon removal, or the use of slowly

186 evolving genes. Recently, however, Salichos and Rokas questioned the exclusive reliance on

187 concatenation and associated practices. In a phylogenomic analysis of 1,070 orthologues from 23

188 yeast genomes, surprisingly, they identified 1,070 distinct gene trees, which were all incongruent

189 with the phylogeny inferred by concatenation [36].

190 As shown in Fig S2, when the trees were drawn using neighbor joining (NJ) methods, the

191 phylogenetic trees of Env inferred from ClustalW, MAFFT, T-coffee, MUSCLE and PRANK

192 alignments are all varied, and CAUSA 2.0 suggested another different evolutionary process.

193 Using MSA aligned by ClustalW, when trees are inferred by maximum likelihood (ML) method

194 in "complete deletion" (CD) mode, as shown in Fig 5A-5F, the DNA tree, the back-translate

195 codon tree and the protein tree are slightly inconsistent; however, when the ML trees were drawn

196 in "using all sites" (AS) mode, the AS trees are more inconsistent with each other, and differ

197 greatly with the CD trees (strain HV1MA was clustered in different clade with a very high

198 Bootstrap value), suggesting that in ClustalW the alignment of the variable region is inconsistent

199 with that of the conserved region; In PRANK, as shown in Fig 5G-5L, the DNA tree, the codon

200 tree and the protein tree are more consistent in CD mode, but the CD trees also differ greatly

201 from the AS trees (HV1MA was classified in different clade with a very high Bootstrap value).

202   In CAUSA 2.0, as shown in Fig 5M-5P, however, the DNA tree and the protein tree are

203   fully consistent in either CD or AS mode, and almost fully consistent between CD and AS mode:

204   the only two inconsistent branches, HV1C4 and HV1A2, stays within a same clade, suggesting

205   that in CAUSA 2.0 the variable region is aligned consistently with the conserved region.

206   In order to further evaluate the accuracy of different alignments and phylogenetic trees in

207   HIV, we build alignments and trees for two genes, env and gag, respectively using ClustalW,

208   PRANK and CAUSA 2.0. As shown in Fig 6A and 6B, phylogenetic trees for gag protein

209   constructed from ClustalW alignment are inconsistent in different mode, but those from

210   PRANK-Codon (Fig 6C, 6D) and CAUSA 2.0 (Fig 6E, 6F) are all fully consistent with each

211   other in different mode. Moreover, the ClustalW protein trees are different between env (Fig 5A)

212   and gag (Fig 6A), suggesting that in some of these HIV genomes, such as HV1JRFL, HV1J3,

213   HV1B1 and HV1C4, the two genes derived from different strains. However, the PRANK CD

214   trees (Fig 5G, 6C) and CAUSA 2.0 CD trees (Fig 5M, 6E) are basically consistent with each

215   other, and suggest that in most of these virus genomes the two different genes underwent a same

216   evolutionary process, except that strain HV1MA might be the only recombinant virus.

217   As described above, CAUSA 2.0 trees are consistent in different tree drawing mode in both

218   of the two HIV genes, Env (Fig 5M-5P) and Gag (Fig 6E, 6F), suggesting that CAUSA 2.0

219   improves the phylogenetic analyses in virus. The trees inferred from PRANK and CAUSA 2.0

220   might be more accurate than the ClustalW trees, because they are more consistent not only in

221   different tree drawing mode, but in different genes. Because of the underlying phylogeny-aware

222   algorithm, PRANK is by far the best alignment method for phylogenetic analysis. CAUSA 2.0

223   might be equivalent to or sometimes even better than PRANK in phylogenetic analysis.

224   *CAUSA 2.0 improves phylogenetic analyses in human and animals*

225   The CAUSA 1.0 algorithm has been applied in POMAGO [37], a multiple genome-wide

226   alignment tool for bacteria. To test the performance of CAUSA 2.0 in human and animals, a set

227   of orthologous proteins families were arbitrarily selected from TreeFam-A, a manually curated

228   database of molecular phylogenetic trees. Firstly, we aligned their CDSs with ClustalW, PRANK

229  and CAUSA 2.0, respectively in DNA, back-translate, codon or unified mode, and then

230  constructed phylogenetic trees using ML method, respectively in "complete deletion" (CD) mode

231  and "using all sites" (AS) mode. Highly conserved protein families without a variable region

232  were excluded, as we observed that different alignments and trees inferred from different

233  methods are basically consistent with each other in those highly conserved proteins. In the

234  variable proteins, which have at least one variable region, we observed that different MSAs and

235  phylogenetic trees differed greatly from each other (Table S3).

236  As shown in Table S3, the topologies of PRANK trees are basically consistent with those of

237  the ClustalW trees, while their bootstrap percentages are better in AS mode. Overall, in these

238  variable protein families tested, CAUSA 2.0 trees are more consistent with the TreeFam

239  reference trees: both the average number of consistent branches (NCBs) and bootstrap

240  percentages (BSPs) are significantly higher than those of the other trees, including protein trees

241  and codon trees. When all sites are used, the CAUSA 2.0-AS trees are the most consistent with

242  the TreeFam reference trees, suggesting that the alignments aligned using CAUSA 2.0 are the

243  most consistent between the variable region and the conserved region.

244  In order to avoid possible biases in human assessment, we further compared the trees with

245  the TreeFam reference trees using TOPD/FMTS [38], a software for the evaluation of similarities

246  between phylogenetic trees. As shown in Table S3, the TOPD/FMTS split distances (SD) are

247  mostly consistent with the number of consistent branches counted. As shown in Table 2, CAUSA

248  2.0 trees shows the lowest SD in AS mode and are significantly better than those of the other

249  alignment methods, implies that CAUSA 2.0 gives more accurate alignments, and brings a

250  higher confidence in the downstream evolutionary analysis of the proteins and their CDSs.

251  ***Testing CAUSA 2.0 on simulated CDSs***

252  As demonstrated in above, CAUSA 2.0 improves the alignment and evolutionary analysis of

253  proteins in real biological data mainly due to the more accurate localization of in-frame indels

254  which is not coincide perfectly with the codon boundaries. In nature, non-in-frame indels, or in-

255  frame indels containing one or more stop codons, are prone to be eliminated in the struggle for

256 existence. Therefore, they might be less frequent in real biological data than predicted. Actually,

257 CAUSA 2.0 simulates such an evolution process by adapting indels to be in-frame as far as

258 possible, thus makes the alignment of coding sequences biologically more significant.

259     We tried to quantify the accuracy of CAUSA 2.0 on simulated CDSs and compare it with

260 those of the other alignment methods. At present, there are only a few programs, such as *indel-*

261 *seq-gen* [39,40], can simulate CDSs with an adjustable indels probability (Id). Unfortunately,

262 however, they does not distinguish in-frame indels from non-in-frame ones. As shown in Fig S3,

263 in the CDSs simulated with low level of indels (Id = 0.01), stop codons were often generated,

264 and all indels coincide perfectly with the codon boundaries; while in the CDSs simulated with

265 high level of indels (Id = 0.05), a lot of non-in-frame indels and frameshift mutations were

266 generated, makes the alignment becomes largely meaningless. In order to test the accuracy of

267 CAUSA 2.0 and the other methods in locating in-frame and non-in-frame indels, a new DNA

268 sequence simulator that can simulate CDSs with a controllable level of in-frame/non-in-frame

269 indels must be developed in future studies.

270 ***Testing CAUSA 2.0 on BAliBASE***

271     We further compared CAUSA 2.0 with the other MSA algorithm on BAliBASE, a hand-

272 curated alignment benchmark [41-43]. As shown in Table S4, the Baliscore for the different protein

273 alignments of BaliBase BB11001 (the high mobility group protein, HMG1) aligned by PRANK,

274 PRANK-codon and CAUSA 2.0 are both lower than those of the other protein alignments. It has

275 been reported that an alignment is not necessarily evolutionarily correct even if it is structurally

276 accurate[7]. Therefore, BAliBASE, which is based on structure alignments, might not be suitable

277 for the assessment of alignment methods aimed at creating evolutionarily correct alignments,

278 such as PRANK, codon alignment or CAUSA 2.0.

279 ***The information contents of different alignments***

280     The information content (IC) for DNA and protein sequences were given by Shannon

281 entropy [44-46]. For unaligned DNA or protein sequences, IC is computed by the frequencies of

282 different states occurred in the sequences (4 bases, 20 amino acids or 64 codons).

283
$$S_B = -\sum_{i=1}^{4} P_i \, log_2(P_i) \qquad [1]$$

284
$$S_A = -\sum_{j=1}^{20} P_j \, log_2(P_j) \qquad [2]$$

285  Given a triplet codon, let $P_k$, $P_l$, $P_m$ be the frequencies of the three bases,

286
$$S_C = -\sum_{k=1}^{4} \sum_{l=1}^{4} \sum_{m=1}^{4} P_k P_l P_m \, log_2(P_k P_l P_m) \qquad [3]$$

287  The information contents of DNA and protein sequences reach their upper limits when

288  different states occur in equal frequencies:

289
$$Max(S_A) = log_2\left(\frac{1}{4}\right) = 2.0000$$

290
$$Max(S_B) = log_2\left(\frac{1}{20}\right) = 4.3219$$

291
$$Max(S_C) = log_2\left(\frac{1}{64}\right) = 6.0000$$

292  After the sequences were aligned, the entropy of the aligned sequences will change, as an

293  additional state (the gap-state) is introduced. Let $S_G$ be the increased IC comes from the gaps

294  inserted during the alignment process, we have:

295
$$S_G = -P_g \, log_2(P_g) \qquad [4]$$

296
$$S_{Base} = S_B + S_G \qquad [5]$$

297
$$S_{AA} = S_A + S_G \qquad [6]$$

298
$$S_{Codon} = S_C + S_G \qquad [7]$$

299  Where $P_i$, $P_j$, $P_k$, $P_g$ is the frequencies of amino acids (A), bases (B), codons (C) and gaps

300  (G) in the aligned sequences, respectively. Given a set of random sequences, in theory the IC of

301  protein sequences is higher than that of DNA sequences, and that of coding sequences is even

302  higher. Therefore, a protein alignment would be better than the DNA alignment, and the

303  corresponding codon alignment could be even better.

304　　　In DNA, protein, codon or back translated protein alignments, there is only one gap-state

305　because codons are indivisible; however, there are many more gap-states in CAUSA 2.0, because

306　every base of a codon can be replaced by a gap. Given a triplet codon, let $P_k$, $P_l$, $P_m$ be the

307　frequencies of the three bases, and $P_g$ be the frequency of a base to be replaced by a gap. Then,

308
$$P_0 = P_k P_l P_m (1 - P_g)^3 \qquad [8]$$

309
$$P_1 = P_k P_l P_m P_g (1 - P_g)^2 \qquad [9]$$

310
$$P_2 = P_k P_l P_m P_g^2 (1 - P_g) \qquad [10]$$

311
$$P_3 = P_k P_l P_m P_g^3 \qquad [11]$$

312　　　Where $P_1$, $P_2$, $P_3$ is the probability of a codon contains 1, 2 and 3 gaps, respectively. The

313　total probability of split codons (the codons contain at least one gaps) is,

314
$$P_s = P_1 + P_2 + P_3 = P_k P_l P_m P_g (1 - P_g + P_g^2) \qquad [12]$$

315　　　The information content of a unified alignment is given by,

316
$$S_u = S_C - \Sigma_{k=1}^4 \sum_{l=1}^4 \Sigma_{m=1}^4 P_s \, log_2(P_s) \qquad [13]$$

317　　　In an unified alignment: (1) If $P_g = 0$, then $P_s = 0$, $S_u = S_{Codon}$, there is no split codon, and

318　thus the information content of the unified alignment equals to that of a codon alignment; In

319　other words, codon alignment could be considered as a special case of the unified alignment,

320　where no split codon is present. (2) If $P_g > 0$, then $P_s > 0$, $S_u > S_{Codon}$, the split codons enhance

321　the information content of the unified alignment, and makes it greater than those of the DNA,

322　protein or codon alignments, thus explains why a unified alignment could be more accurate than

323　the other level of alignments.


324　***The computational efficiency of CAUSA 2.0***


325　　　Among all the programs tested, as shown in Table S5, MAFFT is the fastest, especially

326　when the number of sequences and the total length increases; MUSCLE, an iterative method, is

327　faster than ClustalW. T-Coffee is even faster as it adopted a multiple-threading paralleled

328  computation. PRANK, however, is very time-consuming: even in moderate-sized datasets (20 ~

329  50KB in total length), while other software finished aligning in a few minutes, it takes PRANK

330  several hours to align them. PRANK codon alignment is even more time-consuming, probably

331  due to its adopting of a large (61x61) scoring matrix. In contrast, using a small (24x24) scoring

332  matrix, CAUSA 2.0 is computationally very efficient. The computation time of CAUSA 2.0 is

333  similar to that of ClustalW. However, since CAUSA 2.0 outputs both protein and codon

334  alignment in a single run, it is considered more efficient than most of the other programs.


335  **Discussion and conclusion**


336  Multiple sequence alignment is widely useful in bioinformatics, molecular evolution,

337  genetics and genomics, such as reconstruction of phylogenetic history, structure and functional

338  analyses. Moreover, phylogenetic trees inferred from DNA, protein and codon alignments are

339  often inconsistent. It has been believed that the phylogenetic signal disappears more rapidly from

340  DNA sequences than from their encoded proteins, and therefore it has been preferable to align

341  protein and their CDSs at AA level [8]. However, some important information carried by CDSs,

342  including synonymous substitutions, codon splitting/fusion, siDel and ipFS, get lost when they

343  were translated into AAs, and makes the alignments and phylogenetic trees sometimes inaccurate.

344  CAUSA 2.0 performs a two level dynamic programming alignment algorithm at codon and AA

345  level. CAUSA 2.0 minimizes the total amount of mutations at both codon and AA level

346  efficiently, locates the in-frame indels precisely and better interprets their role in the molecular

347  evolution. In CAUSA 2.0, the boundary of every triplet codon is defined by their encoded amino

348  acid without forcing them into indivisible units. Thanks to the position constraint inherent in the

349  codon-aa 4-tuples, CAUSA 2.0 aligns codons while keeps the corresponding AAs staying in

350  frame without using post-alignment adjustments.

351  In the previous study [27], we have evaluated CAUSA 1.0 for the phylogenetic analysis in

352  virus, bacteria and mammalian animals. However, we used the combined alignments directly for

353  the construction of the phylogenetic trees, and the trees were inferred based on p-distances as

354 input to neighbor joining (NJ) method. It has been pointed out that there are several defects in
355 that study, mainly including: (1) Although CAUSA does improve the alignment of CDSs, there
356 is no model to support using combined alignment directly for the construction of phylogenetic
357 trees; (2) The phylogenetic trees were constructed using Neighbor Joining (NJ) method, but the
358 use of p-distances for phylogenetic analysis has been shown to be systematically biased, because
359 the p-distances are known to undercount the number of substitutions between a pair of sequences,
360 and perform particularly poorly for distantly related sequences; (3) CAUSA is compared only
361 with protein and DNA level alignments, but not with codon level alignments.

362     CAUSA 2.0 is an major update of the CAUSA algorithm, and the present study differs from
363 the previous one mainly in the following aspects: (1) Instead of using a combined alignment
364 directly for the construction of phylogenetic trees, the combined alignment was separated into a
365 protein alignment and a codon alignment, and used to construct the phylogenetic trees, therefore
366 any existing model can still apply in the subsequent evolutionary analysis; (2) The phylogenetic
367 trees were inferred using a Maximum Likelihood (ML) method in the MEGA software; (3)
368 CAUSA 2.0 was compared not only with protein and DNA alignments, but also with codon
369 alignments aligned by PRANK. (4) A mathematical model based on entropy information content
370 was developed, which explains why CAUSA 2.0 works better than DNA, protein and codon
371 level alignments; (5) Several newly discovered mutation mechanisms linked to certain kind of
372 genetic diseases were discussed; (6) It was showed that the computational efficiency of the
373 approach is higher than the other methods.

374     Based on the above analysis, it is concluded that CAUSA 2.0 gives more accurate alignment
375 for proteins and their CDSs, as well as more confident predictions in the phylogenetic and
376 evolutionary analysis. The main benefit of this method is the ability to detect in-frame, as well as
377 non in-frame indels, while keeping the respective amino acids aligned. CAUSA 2.0 is also
378 superior in computational efficiency when compared with the other approaches.

379     Finally, CAUSA 2.0 programs run in Microsoft Windows and Linux, written respectively in
380 Microsoft Visual C# and Java, are available for free in the website www.DNAPlusPro.com.

## Acknowledgements

## References

1    Thompson, J. D., Higgins, D. G. & Gibson, T. J. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic acids research* **22**, 4673-4680 (1994).

2    Larkin, M. A. *et al.* Clustal W and Clustal X version 2.0. *Bioinformatics* **23**, 2947-2948, doi:10.1093/bioinformatics/btm404 (2007).

3    Katoh, K. & Standley, D. M. MAFFT: iterative refinement and additional methods. *Methods in molecular biology* **1079**, 131-146, doi:10.1007/978-1-62703-646-7_8 (2014).

4    Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic acids research* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).

5    Notredame, C., Higgins, D. G. & Heringa, J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. *Journal of molecular biology* **302**, 205-217, doi:10.1006/jmbi.2000.4042 (2000).

6    Loytynoja, A. Phylogeny-aware alignment with PRANK. *Methods in molecular biology* **1079**, 155-170, doi:10.1007/978-1-62703-646-7_10 (2014).

7    Loytynoja, A. & Goldman, N. Phylogeny-aware gap placement prevents errors in sequence alignment and evolutionary analysis. *Science* **320**, 1632-1635, doi:10.1126/science.1158395 (2008).

8    Wernersson, R. & Pedersen, A. G. RevTrans: Multiple alignment of coding DNA from aligned amino acid sequences. *Nucleic acids research* **31**, 3537-3539 (2003).

9    Henikoff, S. & Henikoff, J. G. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences of the United States of America* **89**, 10915-10919 (1992).

10   Styczynski, M. P., Jensen, K. L., Rigoutsos, I. & Stephanopoulos, G. BLOSUM62 miscalculations improve search performance. *Nature biotechnology* **26**, 274-275, doi:10.1038/nbt0308-274 (2008).

11   Gonnet, G. H., Cohen, M. A. & Benner, S. A. Exhaustive matching of the entire protein sequence database. *Science* **256**, 1443-1445 (1992).

12   Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. *Nucleic acids research* **34**, W609-612, doi:10.1093/nar/gkl315 (2006).

13   Abascal, F., Zardoya, R. & Telford, M. J. TranslatorX: multiple alignment of nucleotide sequences guided by amino acid translations. *Nucleic acids research* **38**, W7-13, doi:10.1093/nar/gkq291 (2010).

14   Hein, J. An algorithm combining DNA and protein alignment. *Journal of theoretical biology* **167**, 169-174, doi:10.1006/jtbi.1994.1062 (1994).

15   Hein, J. & Stovlbaek, J. Combined DNA and protein alignment. *Methods in enzymology* **266**, 402-418 (1996).

16   Schneider, A., Cannarozzi, G. M. & Gonnet, G. H. Empirical codon substitution matrix. *BMC bioinformatics* **6**, 134, doi:10.1186/1471-2105-6-134 (2005).

17   Goldman, N. & Yang, Z. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Molecular biology and evolution* **11**, 725-736 (1994).

18   McClellan, D. A. The codon-degeneracy model of molecular evolution. *Journal of molecular evolution* **50**, 131-140 (2000).

19   Zaheri, M., Dib, L. & Salamin, N. A generalized mechanistic codon model. *Molecular biology and evolution* **31**, 2528-2541, doi:10.1093/molbev/msu196 (2014).

20   Miyazawa, S. Superiority of a mechanistic codon substitution model even for protein sequences in phylogenetic analysis. *BMC evolutionary biology* **13**, 257, doi:10.1186/1471-2148-13-257 (2013).

21   McClellan, D. A. The phylogenetic utility of the codon-degeneracy model. *Journal of molecular evolution* **51**, 185-193 (2000).

423  22  Abby, S. S., Tannier, E., Gouy, M. & Daubin, V. Lateral gene transfer as a support for the tree of life. *Proceedings of the National Academy of Sciences*
424  *of the United States of America* **109**, 4962-4967, doi:10.1073/pnas.1116871109 (2012).

425  23  Loytynoja, A. & Goldman, N. An algorithm for progressive multiple alignment of sequences with insertions. *Proceedings of the National Academy of*
426  *Sciences of the United States of America* **102**, 10557-10562, doi:10.1073/pnas.0409137102 (2005).

427  24  Szalkowski, A. M. Fast and robust multiple sequence alignment with phylogeny-aware gap placement. *BMC bioinformatics* **13**, 129, doi:10.1186/1471-
428  2105-13-129 (2012).

429  25  Loytynoja, A., Vilella, A. J. & Goldman, N. Accurate extension of multiple sequence alignments using a phylogeny-aware graph algorithm.
430  *Bioinformatics* **28**, 1684-1691, doi:10.1093/bioinformatics/bts198 (2012).

431  26  Loytynoja, A. & Goldman, N. webPRANK: a phylogeny-aware multiple sequence aligner with interactive alignment browser. *BMC bioinformatics* **11**,
432  579, doi:10.1186/1471-2105-11-579 (2010).

433  27  Wang, X. *et al.* Accurate reconstruction of molecular phylogenies for proteins using codon and amino acid unified sequence alignments (CAUSA).
434  *Nature Precedings <http://hdl.handle.net/10101/npre.2011.6730.1>* (2011).

435  28  Mullaney, J. M., Mills, R. E., Pittard, W. S. & Devine, S. E. Small insertions and deletions (INDELs) in human genomes. *Human molecular genetics* **19**,
436  R131-136, doi:10.1093/hmg/ddq400 (2010).

437  29  Amor, I. M. *et al.* Severe osteogenesis imperfecta caused by a small in-frame deletion in CRTAP. *American journal of medical genetics. Part A* **155A**,
438  2865-2870, doi:10.1002/ajmg.a.34269 (2011).

439  30  Guillaud Bataille, M. *et al.* Systematic screening for PRKAR1A gene rearrangement in Carney complex: identification and functional characterization of
440  a new in-frame deletion. *European journal of endocrinology / European Federation of Endocrine Societies* **170**, 151-160, doi:10.1530/EJE-13-0740
441  (2014).

442  31  Yang, W., Wei, H. & Sang, Y. KCNJ11 in-frame 15-bp deletion leading to glibenclamide-responsive neonatal diabetes mellitus in a Chinese child.
443  *Journal of pediatric endocrinology & metabolism : JPEM* **26**, 743-746 (2013).

444  32  Yang, S. Y. *et al.* EGFR exon 19 in-frame deletion and polymorphisms of DNA repair genes in never-smoking female lung adenocarcinoma patients.
445  *International journal of cancer. Journal international du cancer* **132**, 449-458, doi:10.1002/ijc.27630 (2013).

446  33  Weedon, M. N. *et al.* An in-frame deletion at the polymerase active site of POLD1 causes a multisystem disorder with lipodystrophy. *Nature genetics* **45**,
447  947-950, doi:10.1038/ng.2670 (2013).

448  34  Dudkiewicz, M., Szczepinska, T., Grynberg, M. & Pawlowski, K. A novel protein kinase-like domain in a selenoprotein, widespread in the tree of life.
449  *PloS one* **7**, e32138, doi:10.1371/journal.pone.0032138 (2012).

450  35  Sobhani, M., Tabatabaiefar, M. A., Rajab, A., Kajbafzadeh, A. M. & Noori-Daloii, M. R. Molecular characterization of WFS1 in an Iranian family with
451  Wolfram syndrome reveals a novel frameshift mutation associated with early symptoms. *Gene* **528**, 309-313, doi:10.1016/j.gene.2013.06.040 (2013).

452  36  Salichos, L. & Rokas, A. Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature* **497**, 327-331, doi:10.1038/nature12130
453  (2013).

454  37  Wieseke, N., Lechner, M., Ludwig, M. & Marz, M. in *Bioinformatics Research and Applications*   249-260 (Springer, 2013).

455  38  Puigbo, P., Garcia-Vallve, S. & McInerney, J. O. TOPD/FMTS: a new software to compare phylogenetic trees. *Bioinformatics* **23**, 1556-1558,
456  doi:10.1093/bioinformatics/btm135 (2007).

457  39  Strope, C. L., Abel, K., Scott, S. D. & Moriyama, E. N. Biological sequence simulation for testing complex evolutionary hypotheses: indel-Seq-Gen
458  version 2.0. *Molecular biology and evolution* **26**, 2581-2593, doi:10.1093/molbev/msp174 (2009).

459  40  Strope, C. L., Scott, S. D. & Moriyama, E. N. indel-Seq-Gen: a new protein family simulator incorporating domains, motifs, and indels. *Molecular*
460  *biology and evolution* **24**, 640-649, doi:10.1093/molbev/msl195 (2007).

461  41  Thompson, J. D., Koehl, P., Ripp, R. & Poch, O. BAliBASE 3.0: latest developments of the multiple sequence alignment benchmark. *Proteins* **61**, 127-
462  136, doi:10.1002/prot.20527 (2005).

463  42  Bahr, A., Thompson, J. D., Thierry, J. C. & Poch, O. BAliBASE (Benchmark Alignment dataBASE): enhancements for repeats, transmembrane
464  sequences and circular permutations. *Nucleic acids research* **29**, 323-326 (2001).

465  43  Thompson, J. D., Plewniak, F. & Poch, O. BAliBASE: a benchmark alignment database for the evaluation of multiple alignment programs.
466  *Bioinformatics* **15**, 87-88 (1999).

467  44  Rackovsky, S. & Scheraga, H. A. On the information content of protein sequences. *Journal of biomolecular structure & dynamics* **28**, 593-594;
468  discussion 669-674, doi:10.1080/073911011010524957 (2011).

469  45  Weiss, O., Jimenez-Montano, M. A. & Herzel, H. Information content of protein sequences. *Journal of theoretical biology* **206**, 379-386,
470  doi:10.1006/jtbi.2000.2138 (2000).

471  46  Wu, T. T., Fitch, W. M. & Margoliash, E. The information content of protein amino acid sequences. *Annual review of biochemistry* **43**, 539-566,
472  doi:10.1146/annurev.bi.43.070174.002543 (1974).

473

474

**Table 1**(on next page)

Some disease-related indels with codon splitting/fusion events derived from Human Gene Mutation Database.

| Accession Number | Deletion (^codon number) | Disease | Reference |
|---|---|---|---|
| CD031549 | GCC GTG$^{177}$CGC Aag gca gCC CTG GTC AT | WS | Colosimo (2003)  Hum Mutat 21,622 |
| CD011664 | CTG GTC$^{349}$ATC Ttc tAC CTG TCC TT | WS | Khanim (2001)  Hum Mutat 17,357 |
| CD058021 | TTC GTC$^{416}$ATC Ttc tCC TTC CCC AT | WS | Fang (2005)  Zhonghua Yi Xue Za Zhi 85,2468 |
| CD050485 | CC GCC$^{466}$GGC CTg ctA TCG CTG CTG | WS | Giuliano (2005)  Hum Mutat 25,99 |
| CD021029 | CCC TGC$^{766}$CAC Atc aAG AAG TTC GA | WS | Cryns (2002)  Hum Genet 110,389 |
| CD983013 | GC GTC$^{504}$CCG TGc ctg ctc tat gtc taC CTG CTC TAT | WS | Inoue (1998)  Nat Genet 20,143 |
| CD983491 | GG GCC$^{459}$CTG GCc acc gag gtC ACC GCC GGC | WS | Strom (1998)  Hum Mol Genet 7,2021 |
| CD031551 | GC TTC$^{539}$ATG TGg tgT GAG CTC TCC | WS | Colosimo (2003) Hum Mutat 21,622 |
| CD962054 | GTG GTG$^{166}$TTC Ttcg GGA CGG AGT A | LQTS | Wang (1996)  Nat Genet 12,17 |
| CD076814 | GGC CTC$^{274}$ATC Ttc tCC TCG TAC TT | LQTS | Aizawa (2007)  J CaRdiovasc ElectrOphysiol 18,972 |
| CD044136 | TC ATC$^{275}$TTC TCc tcG TAC TTT GTG | LQTS | Gouas (2004)  Cardiovasc Res 63,60 |
| CD097252 | AC CGA$^{395}$GTA GAa gaC AAG GTA GGC | LQTS | Kapplinger (2009)  Heart Rhythm 6,1297 |
| CD097261 | AC ATG$^{613}$CTT CAc caG CTG CTC TCC | LQTS | Kapplinger (2009) Heart Rhythm 6,1297 |
| CD033996 | GC ATT$^{382}$AAC CCa att gct ctG TAT TTG GTG | HD | Garcia-Barcelo (2003)  Clin Chem  50,93 |
| CI031592 | GAG GAC$^{797}$GAC Gga cTC ACC AAG GA | WS | Colosimo (2003)  Hum Mutat 21,622 |
| CI087227 | AC GCG$^{53}$CCC ATc gcg ccc atC GCG CCC GGC | LQTS | Berge (2008)  Scand J Clin Lab Invest 68,362 |
| CI013349 | GG CAG$^{362}$AAG Cag aag caC TTC AAC CGG | LQTS | Kapplinger (2009)  Heart Rhythm 6,1297 |
| CI983230 | GCC ACG$^{1509}$GCT Tcg gct tCC ATT GAC AT | HI | Nestorowicz (1998)  Hum Mol  Genet 7,1119 |

Note: WS - Wolfram Syndrome; LQTS - Long QT syndrome; HD- Hirschsprung disease; HI – Hyperinsulinism.

**Table 2**(on next page)

Comparing the phylogenetic trees of 30 protein families with their TreeFam reference trees by human eyes and using TOPD/FMTS.

| Software | Alignment Method | Gap or Missing data treatment | Average NCB (Human) | Average SD (TOPD/FMTS) | T-Test Significance |
|---|---|---|---|---|---|
| ClustalW | ClustalW-DNA | CD | 11.37 | 0.1777 | ** |
| | | AS | 11.67 | 0.1498 | * |
| | ClustalW-codon (Back translate) | CD | 11.67 | 0.1512 | * |
| | | AS | 11.83 | 0.1332 | * |
| PRANK | PRANK-DNA | CD | 11.33 | 0.1869 | ** |
| | | AS | 11.37 | 0.1719 | ** |
| | PRANK-codon | CD | 10.47 | 0.3031 | ** |
| | | AS | 11.40 | 0.1865 | ** |
| CAUSA 2.0 | CAUSA 2.0-DNA | CD | 11.47 | 0.1832 | ** |
| | | **AS** | **11.97** | **0.1156** | |

Note:    NCB - Number of Consistent Branches; SD –TOPD/FMTS Split Distance;

CD - Complete deletion; AS - Use all sites; "*" - Significant difference (P<0.05);    "**"

–Extreme significant difference (P<0.01).

**Figure 1**(on next page)

The working flowchart of different strategies for aligning proteins and their coding DNA sequences.

**Figure 2**(on next page)

The protein views of different alignments of a variable region (V2) of Env. HIV or SIV strains were derived from the seed alignment of Pfam gp120 protein family (pf00516).

# Figure 3<sub>(on next page)</sub>

Localization of in-frame indels in different alignment programs. (A, B, C) Codon splitting; (D, E, G) Codon fusion; (G, H, I) siDel and ipFS.

**Figure 4**(on next page)

The CAUSA 2.0 view of an in-frame deletion related to Wolfram Syndrome (CD031549), showing a codon fusion event.

# CD031549 (deletion)

| Seq Name | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----------|-----|-----|-----|------|------|------|------|------|
| Normal | Agcc | Vgtg | Rcgc | Kaag | Agca | Agcc | Lctg | Vgtc |
| Wolfram 1 | Agcc | Vgtg | Rcgc | Ta-- | ---- | --cc | Lctg | Vgtc |
| Wolfram 2 | Agcc | Vgtg | Rcgc | Ta-- | ---- | --cc | Lctg | Vgtc |
| Normal | A | V | R | K | A | A | L | V |
| Wolfram 1 | A | V | R | T | - | - | L | V |
| Wolfram 2 | A | V | R | T | - | - | L | V |
| Normal | gcc | gtg | cgc | aag | gca | gcc | ctg | gtc |
| Wolfram 1 | gcc | gtg | cgc | a-- | --- | -cc | ctg | gtc |
| Wolfram 2 | gcc | gtg | cgc | a-- | --- | -cc | ctg | gtc |

**Figure 5**(on next page)

The maximum likelihood trees for HIV Env protein inferred from different alignments and by different gap/missing data treatment (CD-complete deletion; AS-Use all sites).
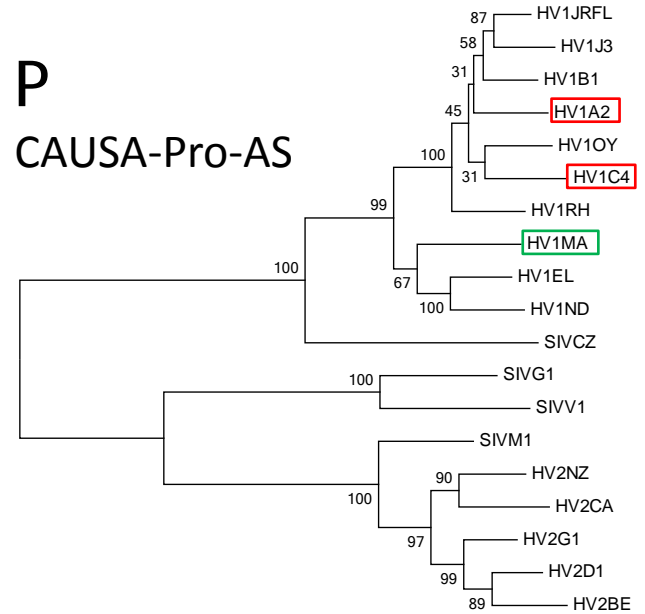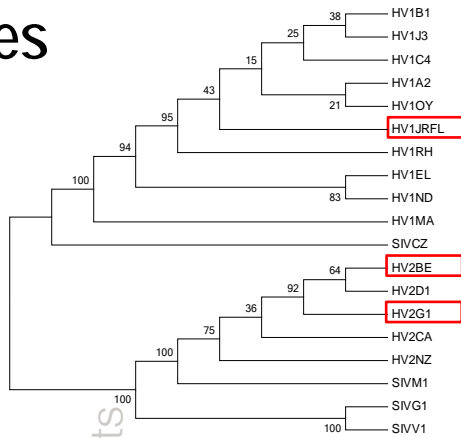
ENV ML trees

M CAUSA-DNA-CD

N CAUSA-Pro-CD

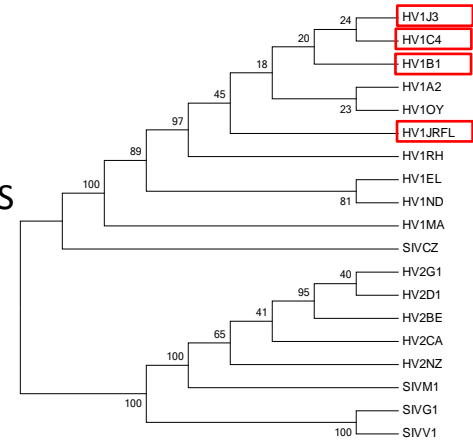O CAUSA-DNA-AS

P CAUSA-Pro-AS

**Figure 6**(on next page)

The maximum likelihood trees for HIV GAG protein inferred from different alignments and by different gap/missing data treatment.
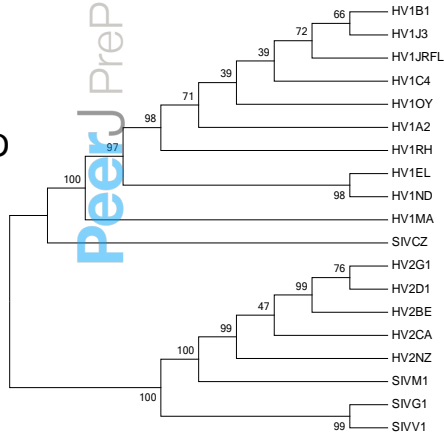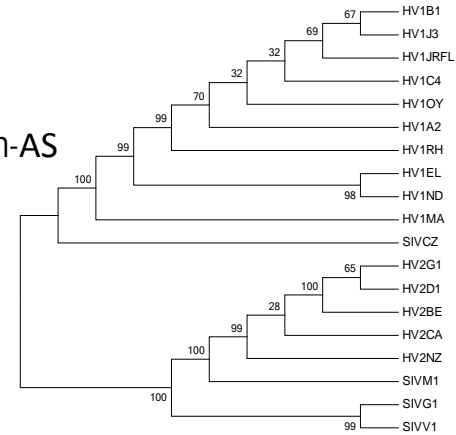
# GAG ML trees
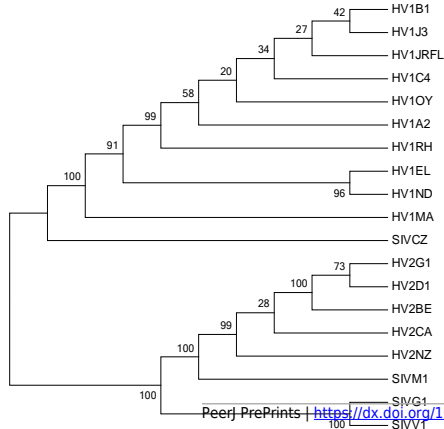


**A** CLUSTALW-**CD**

**B** CLUSTALW-**AS**

**C** PRANK-Codon-**CD**

**D** PRANK-Codon-**AS**

**E** CAUSA-**CD**

**F** CAUSA-**AS**