1    # Curated Multiple Sequence Alignment for the Adenomatous

2    # Polyposis Coli (*APC*) Gene and Accuracy of *In Silico*

3    # Pathogenicity Predictions

4

5    Short Title: Sequence Alignment-Based In Silico Pathogenicity Predictions for APC

6

7    Alexander D. Karabachev[1], Dylan J. Martini[1,#a], David J. Hermel[1#b], Dana Solcz[1], Marcy E.

8    Richardson[2], Tina Pesaran[2], Indra Neil Sarkar[3,4], Marc S. Greenblatt[1]*

9

10   [1] Department of Medicine, University of Vermont, Larner College of Medicine, Burlington VT;

11   [2] Ambry Genetics, Aliso Viejo, CA;

12   [3] Center for Biomedical Informatics, Brown University, Providence, RI

13   [4] Rhode Island Quality Institute, Providence, RI

14   [#a] Current Address: Emory University School of Medicine, Atlanta, GA;

15   [#b] Current Address: Keck School of Medicine of USC, Los Angeles, CA;

16

17   *Corresponding Author:

18   Email: Marc.Greenblatt@uvmhealth.org  (MSG)

19

20

21

22   **Abstract**

23   Computational algorithms are often used to assess pathogenicity of Variants of Uncertain

24   Significance (VUS) that are found in disease-associated genes.  Most computational methods

25   include analysis of protein multiple sequence alignments (PMSA), assessing interspecies

26   variation.  Careful validation of PMSA-based methods has been done for relatively few genes,

27   partially because creation of curated PMSAs is labor-intensive. We assessed how PMSA-based

28   computational tools predict the effects of the missense changes in the *APC* gene, in which

29   pathogenic variants cause Familial Adenomatous Polyposis.  Most Pathogenic or Likely

30   Pathogenic APC variants are protein-truncating changes. However, public databases now

31   contain thousands of variants reported as missense. We created a curated APC PMSA that

32   contained >3 substitutions/site, which is large enough for statistically robust *in silico* analysis.

33   The creation of the PMSA was not easily automated, requiring significant querying and

34   computational analysis of protein and genome sequences.  Of 1924 missense APC variants in

35   the NCBI ClinVar database, 1800 (93.5%) are reported as VUS. All but two missense variants

36   listed as P/LP occur at canonical splice or Exonic Splice Enhancer sites. Pathogenicity

37   predictions by five computational tools (Align-GVGD, SIFT, PolyPhen2, MAPP, REVEL) differed

38   widely in their predictions of Pathogenic/Likely Pathogenic (range 17.5–75.0%) and

39   Benign/Likely Benign (range 25.0–82.5%) for *APC* missense variants in ClinVar. When applied

40   to 21 missense variants reported in ClinVar as Benign, the five methods ranged in accuracy

41   from 76.2-100%. Computational PMSA-based methods can be an excellent classifier for

42   variants of some hereditary cancer genes. However, there may be characteristics of the *APC*

43   gene and protein that confound the results of *in silico* algorithms. A systematic study of these

44   features could greatly improve the automation of alignment-based techniques and the use of

45   predictive algorithms in hereditary cancer genes.

46

47    Author Summary

48    A critical problem in clinical genetics today is interpreting whether a genetic variant is benign or

49    causes disease (pathogenic).  Some of the hardest variants to interpret are those that change

50    one amino acid for another in a protein sequence (a "missense variant"). Various computer

51    programs are often used to predict whether mutations in disease-associated genes likely cause

52    disease.  Most computer programs involve studying how the gene has changed during

53    evolution, comparing the protein sequences of different species by aligning them with each

54    other.  Variants in amino acids that have not tolerated mutation during evolution are usually

55    predicted to be pathogenic, and variants in amino acids that have tolerated variation are usually

56    predicted to be benign. High quality alignments are necessary to make accurate predictions.

57    However, creating high quality alignments is difficult, not easily automated, and requires

58    significant manual curation. Results from computer-generated predictions are used in current

59    published guidelines as one tool for evaluating whether variants will disrupt the protein function

60    and cause disease. These guidelines may be applied to genes in which single amino acid

61    substitutions do not commonly cause disease.  One such example is the APC gene, which is

62    responsible for Familial Adenomatous Polyposis (FAP).  Missense APC changes are not a

63    common cause of FAP. Our analysis of APC demonstrated the difficulty of generating an

64    accurate protein sequence alignment and the tendency of computer tools to overestimate the

65    damaging effects of amino acid substitutions.  Our results suggest that the rules for using

66    computer-based tools to predict whether a variant causes disease should be modified when

67    applied to genes in which missense variants rarely cause disease.

68    **Introduction**

69    Multi-gene panel testing is now routine for identifying hereditary cancer susceptibility, leading to

70    increased detection of pathogenic mutations, which can improve clinical management.

71    However, testing often identifies variants of uncertain significance (VUS), which are often

72    missense amino acid (AA) substitutions, small in frame deletions and duplications, or non-

73    coding changes [1, 2].  VUS in genes that predispose to hereditary cancer and other disorders

74    are rapidly accumulating in variant databases.  For example, the ClinVar database at the

75    National Center for Biotechnology Information at the United States National Library of Medicine

76    provides a freely accessible archive of variants with assertions regarding the pathogenicity of

77    each variant with the indicated phenotype from submitting laboratories and expert panels [3].

78    The classification of these VUS represents a major challenge in clinical genetics.

79    Computational *(in silico)* tools have been developed to help predict whether or not the protein

80    function will be disrupted (reviewed in [5]).  *In silico* tools often use Protein Multiple Sequence

81    Alignments (PMSA) to consider the evolutionary conservation and biophysical properties of the

82    wild type and variant protein to make predictions of pathogenicity. PMSA-based computational

83    methods are complicated to use properly (reviewed in [5]).  The PMSA must be of high quality

84    and sample enough species to provide reliable data [6,7].  These *in silico* methods have been

85    validated for relatively few hereditary cancer genes in which pathogenic missense variants are

86    not rare (BRCA1/2, the mismatch repair [MMR] genes, TP53, a few others) [6, 8, 9, 10, 11, 12].

87    They have not often been validated for other genes, and for some genes predictive value was

88    not strong [13].  However, they are often cited as evidence in favor or against pathogenicity of

89    variants for genes in which validation is lacking.  The American College of Genetics and

90    Genomics (ACMG) and the Association for Molecular Pathology (AMP) published guidelines for

91    evaluating the pathogenicity of variants in Mendelian disease genes, including general rules for

92    the use of *in silico* tools [4].

93    Missense pathogenic variants are rare in some genes, including *APC*, the gene responsible for

94    Familial Adenomatous Polyposis (FAP).  *APC* has been sequenced frequently in clinical genetic

95    testing, but few missense pathogenic variants have been identified, for reasons that have not

96    been clearly demonstrated [14].  The increase in clinical DNA sequencing tests for cancer

97    predisposition has led to an increase in missense VUSs in APC that require classification.

98    Here we systematically apply *in silico* methods to *APC*, assessing the logistics and results of

99    using these commonly available tools to predict pathogenicity of missense variants in a gene for

100   which missense is an uncommon mechanism of pathogenicity.

101

102   **RESULTS**

103   **PMSA Creation**

104   Results from searching the NCBI Gene database for "APC" initially yielded reliable full length

105   APC protein sequences from 38 organisms.  We encountered a number of challenges to the

106   simple automated assembly of a meaningful APC PMSA, including:

107   a)  *Large inconsistencies with the APC human sequence*.  In order to include only sequences

108   which accurately reflect human biology, such sequences were omitted.

109   b)  *Multiple APC isoforms were found for 21 organisms*. To choose the most appropriate

110   isoform, all 104 sequences were aligned using Clustal2W. Isoforms that lacked a common

111   beginning protein sequence of MAA were deleted (N=26). When duplicate sequences were

112   found for the same species, the more complete sequence was used, and if similar length

113   isoforms of the same organism were found with a common sequence initiation, the lowest

114   number isoform was chosen.

115   c)  *Large deletions or insertions*.  Many of these could easily be identified as errors in automated

116   identification of exon-intron boundaries.  In most cases we could identify the appropriate

117   boundary and either insert or delete the appropriate sequence.  For insertions that were unique

118    to one organism, especially in areas of otherwise high homology, BLAST was used to seek

119    other homologues of the inserted sequence, and assessed the relevant nucleotide sequence for

120    plausible overlooked splice sites.

121    d) *Small deletions or insertions.* Short gaps that were confirmed to occur distant from an exon-

122    intron boundary were allowed. The longest such gap was AA 1631-1637 in *Loxodonta africana*

123    (African elephant) and *Trichechus manatus latirostris* (Florida manatee), a highly conserved

124    region in other sequences. Because of the close taxonomic relationship between these two

125    organisms, and the fact that their sequence was assembled on the same Broad Institute

126    platform as many other species in our alignment that lack the deletion, we assessed this gap as

127    likely real.

128    We constructed two PMSAs. Our goal was to create a curated PMSA that would optimize

129    predictions for pathogenicity of variants from computational algorithms. This 10-sequence

130    PMSA contained species chosen to reflect as closely as possible the 14-species PMSA

131    previously reported for analyzing variants and validating computational algorithms in the MMR

132    genes, in which missense VUS are common and *in silico* interpretation is frequently used [8].

133    We identified full length APC sequences for 11 of these 14 species. The 10-species PMSA that

134    we curated using the above criteria (Table 1, PMSA excerpt in Figure 1, full PMSA in

135    Supplementary Figure 1) contained five mammalian APC sequences plus chicken (*Gallus*

136    *gallus*), frog (*Xenopus laevis*), zebrafish (*Danio rerio*), sea urchin (*Strongylocentrotus*

137    *purpuratus*), and sea squirt (*Ciona intestinalis*). A larger PMSA with the full set of 38 full length

138    sequences also was constructed, with reconstitution of obvious missing exons but no detailed

139    curation (Supplementary Figure 2).

140

141

142    **Table 1**. APC amino acid sequences from the NCBI database used in the ten species APC

143    Protein Multiple Sequence Alignment (PMSA) and phylogenetic tree.

| Species | APC |
|---|---|
| Human (*Homo sapiens)* | AAA03586.1 |
| Monkey (*Macaca mulatta*) | XP_014996065.1 |
| Cow (*Bos taurus*) | NP_001069454.2 |
| Mouse (*Mus musculus*) | NP_031488.2 |
| Opossum (*Monodelphis domestica*) | XP_007497871.1 |
| Chicken (*Gallus gallus*) | XP_004949340.1 |
| Frog (*Xenopus laevis*) | NP_001084351.1 |
| Zebrafish (*Danio rerio*) | NP_001137312.1 |
| Sea urchin (*Strongylocentrotus purpuratus*) | XP_783363.3 |
| Sea squirt (*Ciona intestinalis*) | XP_018668496.1 |

144

145    Manual curation was often necessary to identify and label correct exon-intron boundaries and

146    address insertions, gaps, and poorly-conserved areas where the alignment was less certain.  A

147    small amount of manual curation of gaps and insertions was required for vertebrate species.

148    The intronic regions flanking large insertions were examined and assessed as potential splice

149    sites. Sites with a high splice score (see Methods) were interpreted as actual splice sites and

150    retained for creation of the phylogenetic tree. Inserted sequences flanked by a lower than

151    average splice site were omitted from further analyses.

152    More extensive manual curation was required for *C. intestinalis* and *S. purpuratus,* the most

153    distant species used, to ensure an accurate alignment and tree.  Using BLAST+ on insertions in

154    sea squirt and sea urchin that were not present in the human sequence, we identified

155    sequences with little homology on inspection to the vertebrate APC sequences.  Exon 1 (M1 to

156    Q46) and Exon 5, 6, and 7 (A265 to K414) of sea squirt (*C. intestinalis*) and exon 6 (A260 to

157    F477) of sea urchin *(S. purpuratus*) did not align with the other APC sequences, returned

158    negative BLAST results, and were removed from the final PMSA. A region of of *S.purpuratus*

159    was found with homology to a spindle fiber sequence, and a long region in its C-terminus was

160    homologous to a herpesvirus sequence.  Because the exons containing these sequences also

161    contained regions with high homology to APC, the full exons were retained in our PMSA.  A

162    large insertion in *S. purpuratus* containing many consecutive glutamines presumably represents

163    a coding region microsatellite. Sequences flanking this insertion were found with high splice

164    scores, so it was kept in the alignment.

165

166    **Evolutionary rate of APC:**

167    To predict if a given invariant position is invariant with statistical significance (>95% probability),

168    the PMSA must contain >3.0 substitutions/site [6, 7]. In addition to our ten-sequence PMSA,

169    curated alignments were created of nine and eight sequences that omitted the more distant

170    species *Ciona intestinalis* (sea squirt) and *Strongylocentrotus purpuratus* (sea urchin) (data not

171    shown). Applying the PHYLIP ProtPars package to the curated 8, 9, and 10 species APC

172    PMSAs, we calculated that our ten species curated *APC* alignment contained 3.3 substitutions

173    per site (subs/site), sufficient for proceeding with subsequent analyses (see Methods).  Both

174    eight- and nine-sequence PMSAs, omitting the nonvertebrate species, contained fewer than

175    three subs/site.  We calculated subs/site for six other PMSAs of cancer susceptibility genes

176    using the same 10 species found on the Align-GVGD website (Table 2).  *APC* had a

177    comparable evolutionary rate with *CHEK2* and *PMS2*, whereas three MMR genes (*MLH1,*

178    *MSH2, MSH6*) were better conserved (1.6-2.1 subs/site), and *RAD51* was the most well-

179    conserved of the seven genes (0.62 subs/site).

180

181    **Table 2**. Substitutions per site in PMSAs of seven hereditary cancer genes using 10 species

182    with evolutionary depth to sea squirt calculated using the PHYLIP ProtPars package.

183

| Protein | Substitutions per site |
|---------|------------------------|
| PMS2    | 3.4                    |
| APC     | 3.3                    |
| CHEK2   | 3.2                    |
| MSH6    | 2.8                    |

| | |
|---|---|
| MLH1 | 2.1 |
| MSH2 | 1.6 |
| RAD51 | 0.62 |

184

**Phylogenetic Tree Construction**

186   Phylogenetic trees were generated using Bayesian, Maximum Likelihood, and Maximum

187   Parsimony -based methods.  The methods yielded similar trees, and the Maximum Parsimony -

188   based examples are displayed in Figures 2A (10 species) and 2B (38 species). The

189   relationships of the *APC* sequences among different species was as expected with sea urchin

190   and sea squirt as the most distantly related organisms to humans.

191

*APC* **Variants from Public Databases**

193   In the LOVD database maintained by the International Society for Gastrointestinal Hereditary

194   Tumors (InSiGHT), in July 2013 there were a total of 46 *APC* missense variants.  In ClinVar in

195   July 2018, there were a total of 4891 *APC* variants, of which 1988 are missense. Using filters of

196   "missense, pathogenic, likely pathogenic", yielded nine variants in the ClinVar database with

197   assertions of Pathogenic/Likely Pathogenic (P/LP) and no conflicting interpretations of

198   pathogenicity per ClinVar criteria.  Upon further examination, it was determined that two variants

199   were somatic mutations, and the pathogenicity of the other seven variants were inferred to be

200   from a splicing abnormality.  Six were found to occur at canonical splice sites, and the seventh

201   occurs within an Exonic Splicing Enhancer sequence, with confirming RNA and *in vitro* evidence

202   of splicing alterations [15] (Supplementary Table).  Thus, no pathogenic missense germline

203   *APC* variants were documented in ClinVar using these search parameters.  There are n=21

204   variants (1.3% of all missense variants) with assertions of Benign or Likely Benign (B/LB).  All of

205   these were classified using criteria other than *in silico* algorithms.  Of the remaining variants in

206    ClinVar, 93.5% of the missense variants are reported as "Unknown Significance"; the rest are

207    classified as either "Other", or display conflicting assertions of pathogenicity (Table 3).

208

209    **Supplementary Table 1.** Nine APC missense variants using filters for "missense, pathogenic,

210    likely pathogenic".

| APC Classified Pathogenic Variant | ClinVar Classification | Type of Variant |
|---|---|---|
| R141S | Pathogenic | Splice Site |
| K516N | Pathogenic | Splice Site |
| K581N | Likely Pathogenic | Splice Site |
| S634R | Likely Pathogenic | Exonic Splice Enhancer site |
| R653M | Pathogenic | Splice Site |
| R653G | Pathogenic | Splice Site |
| R653K | Pathogenic | Splice Site |
| G1120E | Pathogenic | Somatic |
| S1395C | Pathogenic | Somatic |

211

212    Supplementary Table 1 Legend: Two are due to somatic mutations, six are located in canonical

213    splice sites and one occurs within an Exonic Splicing Enhancer sequence. [p.S1028N was

214    reclassified as pathogenic in the ClinVar database after we collected the data – Not included]

215    https://preview.ncbi.nlm.nih.gov/clinvar/variation/428186/]

216

217    **Table 3.** APC missense variants from the NCBI ClinVar database with Clinical Significance

218    Classifications of: "Benign", "Likely Benign", "Pathogenic", "Likely Pathogenic", "Uncertain

219    Significance" and "Conflicting Interpretations of Pathogenicity".

220

| ClinVar "Clinical Significance" for APC | Missense Variants (N=1924) |
|---|---|
| Benign/Likely Benign | 21 (1.1%) |
| Pathogenic/Likely Pathogenic | 0 (0%) |
| Uncertain Significance | 1800 (93.5%) |
| Conflicting Interpretations of Pathogenicity | 103 (5.4%) |

221

222    Table 3 Legend: Substitutions flanking the 12 splice sites found in Human APC were removed

223    from the list of selected missense variants. A total of 1924 variants that met the above

224    classification criteria and were not located in exon boundaries were used for analysis. Of the

225    1924 variants, 1.1% were classified as benign, none were classified as pathogenic and 98.9 %

226    were classified as uncertain or conflicting interpretation of pathogenicity.

227

228    **Computational methods to classify APC variants**

229    To predict the pathogenic effects of missense substitutions, multiple computational algorithms

230    based on PMSAs and evolutionary conservation have been developed. We applied five of these

231    tools (SIFT, PolyPhen2, Align-GVGD, MAPP, REVEL) to analyze *APC* missense variants.

232    For the n=21 variants classified in ClinVar as B/LB, the prediction algorithms showed good

233    concordance with each other and with the ClinVar classifications (Table 4A).  REVEL and A-

234    GVGD showed 100% concordance with ClinVar, SIFT predicted 95.5%, PolyPhen2 81.8%, and

235    MAPP 77.8% to be Neutral.  For the n=1904 variants classified as VUS, "Other", or conflicting,

236    the output differed significantly among the four non-aggregating methods (excluding REVEL).

237    The proportion of variants predicted to be "Benign" were MAPP 25.0%, PolyPhen2 41.0%, SIFT

238    68.1%, Align-GVGD 82.5% (Table 4A). For MAPP, we initially used the cutoff score of 4.5

239    previously established to distinguish P/LP from B/LB *MLH1* and *MSH2* variants [8]. This cutoff

240    predicted 75% of APC VUS to be pathogenic, an improbable proportion. With no known

241    pathogenic missense variants, it is unclear what cutoff score is appropriate. The lowest MAPP

242    cutoff score (34.79) that achieved a specificity and total accuracy of 100% for classifying benign

243    variants predicts 2.6% of VUS as pathogenic.

244

245 **Table 4A. Predictions of substitution severity with different *in silico* programs**

246

| | | Benign Variants (N=21) | | | VUS (N=1904) |
|---|---|---|---|---|---|
| Method | Classification | Total (%) | Specificity | Total Accuracy | Predictions: Total (%) |
| ClinVar | Pathogenic | 0 (0%) | - | - | - |
| | Benign | 21 (100%) | | | - |
| REVEL | Deleterious (REVEL score ≥ 0.5) | 0 (0%) | 100% | 100% | N/A |
| | Neutral (REVEL score < 0.5) | 21 (100%) | | | N/A |
| A-GVGD | Class C65 (Deleterious moderate) | 0 (0%) | 100% | 100% | **77 (4.0%)** |
| | Class C55 (Deleterious supporting) | 0 (0%) | | | **37 (1.9%)** |
| | Class C45 (Deleterious supporting) | 0 (0%) | | | **8 (0.42%)** |
| | Class C35 (Deleterious supporting) | 0 (0%) | | | **27 (1.4%)** |
| | Class C25 (Deleterious supporting) | 0 (0%) | | | **64 (3.3%)** |
| | Class C15 (Deleterious supporting) | 0 (0%) | | | **120 (6.3%)** |
| | Class C0 (Neutral) | 21 (100%) | | | **1571 (82.5%)** |
| SIFT | Deleterious | 1 (4.8%) | 95.4% | 95.2% | **608 (31.9%)** |
| | Tolerated | 20 (95.2%) | | | **1296 (68.1%)** |
| PolyPhen2 | Probably Damaging | 1 (4.8%) | 84.0% | 80.9% | **814 (42.8%)** |
| | Possibly Damaging | 3 (13.3%) | | | **309 (16.2%)** |
| | Benign | 17 (80.9%) | | | **781 (41.0%)** |
| MAPP | Pathogenic (MAPP score ≥ 4.5) | 5 (23.8%) | 80.7% | 76.2% | **1428 (75.0%)** |
| | Neutral (MAPP score < 4.5) | 16 (76.2%) | | | **476 (25.0%)** |

247

248 Table 4A Legend: Predictions of pathogenicity for APC missense variants were made using

249 REVEL, A-GVGD, SIFT, PolyPhen2 and MAPP. REVEL output classes were designated as

250 "Deleterious" for variants with a REVEL score ≥ 0.5 and "Neutral" with a REVEL score < 0.5

251 [16]. Assigning A-GVGD output Classes as "Neutral", "Deleterious moderate" and "Deleterious

252 supporting" are based on probabilities from [17] and quantitative modeling of the ACMG/AMP

253 criteria for assigning pathogenicity [4, 18]. SIFT predicts substitutions with SIFT scores less than

254 0.05 as "Deleterious" and scores equal to or greater than 0.05 as "Tolerated" [19]. PolyPhen2

255 predicts variants based on a Position Specific Independent Count (PSIC) score as "Benign" and

256 "Probably Damaging" with high confidence, while a prediction of "Possibly Damaging" is

257 predicted to be damaging, but with low confidence [20]. For MAPP, we used a cutoff score of

258 4.5 to predict "Pathogenic" versus "Neutral" substitutions based the cutoff used to distinguish

259 pathogenic and neutral variants for MLH1 and MSH2. [8].

260

261    We explored the hypothesis that protein structural features would be associated with the

262    likelihood that a VUS was pathogenic or benign.  APC contains multiple repeats of the β-catenin

263    binding and armadillo repeats, plus domains for oligomerization, and binding to microtubules,

264    and EB1 and DLG proteins [21].  We hypothesized that missense variants 1) in the β-catenin

265    binding and armadillo repeats would be neutral, since there was domain redundancy, 2) in the

266    non-repeated domains would be more likely to be pathogenic, and 3) in unstructured regions

267    would be neutral.  There was no difference in the distribution of variants classified in ClinVar as

268    neutral versus VUS relative to the beta catenin, armadillo, or other domains (Table 4b).

269

270    **Table 4B**. Proportion of Benign/Likely Benign variants and Variants of Unknown Significance by
271    APC Protein Structural Feature.

| Domain | Benign/Likely Benign | Unknown Significance |
|---|---|---|
| Beta catenin | 5 (23.8%) | 606 (31.8%) |
| Armadillo | 1 (4.8%) | 156 (8.2%) |
| Other domains | 4 (19.0%) | 378 (19.9%) |
| Not in domain | 11 (52.4%) | 764 (40.1%) |
| Total | 21 | 1904 |

272

273    Per our examination of the ClinVar database in May 2018, all APC missense mutations noted as

274    P/LP were found to be somatic mutations, or located in canonical splice sites, or located in

275    Exonic Splicing Enhancer sequences.  Shortly after we closed our data set, p.S1028N, located

276    in the first of four highly conserved 15-amino acid repeats within the β-catenin binding domain,

277    was submitted to ClinVar by Ambry Genetics and classified as Likely Pathogenic. The evidence

278    for this classification includes, as per the ACMG/AMP guidelines, segregation score

279    (PP1_Strong, six meioses), phenotype score (PS4_Moderate), functional domain (PM1 [22]),

280    population frequency score (PM2_Supporting) and *in silico* data (PP3).  There is no evidence of

281    splice abnormality.  This variant would reach LP regardless of *in silico* analysis.  Further scrutiny

282    of variants in this region demonstrates one other variant, p.N1026S, classified as "Conflicting

283    Interpretations of Pathogenicity" in ClinVar, which satisfies the ACMG/AMP guidelines as LP.

284    The same criteria (PP1_Strong, PS4_Moderate, PM1, PM2) can be applied to p.N1026S, in

285    addition to a functional defect (PS3) as reported in the literature [23, 22]. N1026 and S1028 are

286    both located in the first 15-amino acid repeat of the β-catenin binding domain and after careful

287    review are the only LP/P *APC* missense variants that we found in ClinVar in July 2018 that

288    satisfy the ACMG/AMP guidelines.

289

290    **DISCUSSION**

291    *In silico* tools have been validated with accepted standards for relatively few genes, and the field

292    would greatly benefit from refinement of standards for applying these tools.  Factors that have

293    been shown to be important for interpreting the output and reliability of computational algorithms

294    include quality of PMSA (reviewed in [5]), and choice of variant data sets [24].  An important

295    factor regarding data sets that has emerged recently is how predictors should not be evaluated

296    on variants or proteins that were used to train their prediction models.  This circularity could

297    result in predictive values that are artificially inflated [24, 25], and could occur with either likely

298    pathogenic or likely benign variants. We suggest that not enough attention has been assigned

299    to an additional important factor, the likelihood that missense substitution is a major mechanism

300    of pathogenicity for a gene.

301    Our analysis suggests possible revisions to the ACMG/AMP classification scheme for

302    pathogenicity, which defines multiple criteria for evidence of benign or pathogenic effect, with

303    strength ranging from "Supporting" to "Very Strong", and rules for combining different types of

304    evidence [4].  For example, criterion BP1, "Missense variant in a gene for which primarily

305    truncating variants are known to cause disease", is relevant to *APC*.  By this criterion, any

306    missense APC variant is given "Supporting" evidence, the lowest level, favoring benign

307    classification of missense variants.  Further study may help determine whether this criterion for

308    benign classification should be upgraded from "Supporting" (for which estimated Odds of

309    Pathogenicity is low [18], discussed below) to a higher level for these variants.  The PP2

310    criterion for pathogenicity presupposes that missense is a common mechanism for mutation;

311    future studies should assess whether it is being inappropriately used when missense is a rare or

312    unknown mechanism for a given gene.

313    Our work confirms that PMSA construction remains a labor-intensive task [26].  Current

314    automated tools do not align unstructured regions accurately, resulting in errors that require

315    manual curation of protein and nucleotide sequences in order to optimally curate a full

316    alignment.  For many genes, accurate PMSA can prove important for *in silico* analysis of variant

317    pathogenicity [5].  There is no consensus in the assessment of PMSA quality, although metrics

318    have been proposed [27].  We and others have proposed that a PMSA should include enough

319    sequences to contain three subs/site in order for predictions to be statistically robust [6, 7], and

320    for APC we achieved this threshold with the addition of non-vertebrate sequences.  We chose

321    our sequences to be consistent with PMSAs of other cancer susceptibility genes for which *in*

322    *silico* algorithms have proven to be valuable tools for variant classification. PMSAs for 15 such

323    genes are posted on the Align-GVGD (http://agvgd.hci.utah.edu/about.php) web site.  We hope

324    to promote standardization of methods for the purposes of *in silico* analysis for variant

325    classification.  It remains to be determined whether a consistent set of sequences will be most

326    appropriate for other gene sets. The creation and validation of our APC PMSA did identify

327    interesting features of gene evolution and of genome annotation and analysis, and we anticipate

328    that PMSAs across gene families are likely to elucidate specific structure-function relationships

329    and molecular pathways of critical cellular functions. The full APC PMSA can be seen in

330    Supplementary Figure 1, where it can be used for purposes that are beyond the scope of this

331    paper.

332     One cannot assume that *in silico* tools that are valuable predictors for one gene will perform as

333     well for other genes.  The majority of *APC* missense variants in ClinVar are likely to be benign,

334     given the paucity of missense pathogenic variants identified in over two decades of clinical *APC*

335     testing.  An example of a similar gene is *CDH1*, in which pathogenic missense variants also are

336     rare.  An expert panel studying the *CDH1* gene has recommended that computational methods

337     not be used for missense *CDH1* variants [28].  Thus, tools that work well for genes that are

338     commonly inactivated by missense changes [29, 30, 8, 12] can be misleading for genes that are

339     rarely inactivated by missense.  For such genes, traditional *in silico* tools will likely overestimate

340     the probability of pathogenicity of any missense variant.

341     The ClinGen Sequence Variant Interpretation working group has estimated that the "Supporting"

342     level of evidence confers approximately 2.08/1 odds in favor of pathogenicity [18], or a 67.5%

343     probability of pathogenicity.  Our current analyses of APC variants suggest that the likelihood

344     that a missense APC variant is pathogenic is far lower than 1%.  Despite this, our curated APC

345     PMSA and several *in silico* prediction tools all predicted a significant fraction of missense

346     variants to be pathogenic. The methods that we used varied widely in their predictions for APC

347     VUS; predictions of Pathogenic or Likely Pathogenic ranged from 17.5% to 75%, all of which are

348     higher than the likely figure by at least an order of magnitude.  This provides mathematical

349     support for not using *in silico* evidence in favor of pathogenicity (PP3 in the ACMG/AMP scheme

350     [4]) for these genes.  One approach might be to create a decision tree in which a gene must

351     meet specific criteria before *in silico* evidence is applied.  More work is needed in order to

352     understand which genes require pre-curation to assess whether PMSA-based or other *in silico*

353     methods are likely to be useful.  A difference between functional or structural relevance to the

354     protein and clinical relevance may occur if the assayed function is not crucial to the phenotype,

355     or perhaps from domain redundancy or other protein structural features.

356    Another important factor regarding data sets is whether the subject was being tested because of

357    clinical suspicion, or whether broad panel testing, whole exome or whole genome sequencing

358    yielded a variant in the absence of any known clinical features.  The degree of clinical suspicion

359    is difficult to discern from the majority of ClinVar *APC* variants.  The prior probability of

360    pathogenicity [8] will be much lower for a variant discovered incidentally through whole exome

361    sequencing compared with one identified through clinical testing because of a strong history of

362    polyposis and/or colon cancer, with intermediate scenarios also possible.

363    Computational methods can be an excellent classifier for missense variants in hereditary cancer

364    genes where missense is a common mechanism of pathogenicity [8-12].  However, known

365    pathogenic APC missense germline variants are rare.  It is possible that none exist outside of

366    the first 15-amino acid repeat of the β-catenin binding domain, and it is unknown how many

367    other pathogenic missense variants are located in this 15 amino acid repeat, complicating the

368    use of computational tools.  Further analysis of this region is necessary to understand the role of

369    missense APC variants and the value of *in silico* algorithms. The β-catenin binding repeats may

370    be the only specific region of 15 AA out of the 2843 AA of APC in which *in silico* methods may

371    be predictive of clinical pathogenicity. A similar observation to the use of *in silico* analysis has

372    been made regarding the BRCT domain of BRCA1 [5].  There may be characteristics of the

373    APC gene and protein that confound the results of *in silico* algorithms. One plausible hypothesis

374    for the failure of missense variants to abrogate APC function is the redundancy of APC

375    important structural elements (armadillo repeats, β-catenin and axin binding sites) [21], so the

376    inactivation of a single repeat might not eliminate binding to the target to a clinically relevant

377    level.

378    Defining features that distinguish genes for which missense is a common (e.g., MMR genes [8])

379    versus uncommon (e.g., *CDH1* [31, 28], *RB1* [32]) pathogenic mechanism would significantly

380    improve the application of *in silico* tools to variant classification.  We propose that *in silico*

381    methods to assess missense variants (PP3 and BP4 in the ACMG/AMP guidelines [4]) be used

382    sparingly for any gene where strong evidence suggests that missense rarely causes

383    pathogenicity.  Future work might consider whether BP4 (concordance for "benign" classification

384    among multiple methods) might be replaced by BP1 (truncation predominates, missense

385    unlikely) in such cases.  Our results suggest that a systematic study of variant pathogenicity and

386    protein features such as domain structure is warranted to improve the use of predictive

387    algorithms in hereditary cancer genes.

388

389

390    **Methods**

391    Sequence and variant data are publically available from databases at the NLM.  The study

392    protocol was determined to be exempt from human subject regulations by Western IRB, as the

393    data were de-identified.

394    **APC Sequences and Multiple Sequence Alignments, Phylogenetic analysis**

395    Amino acid sequences were collected by searching NCBI's online Gene database

396    (http://www.ncbi.nlm.nih.gov/gene), for "APC" in 2013, 2015, and 2018.  PMSAs were made

397    using Clustal Omega from the European Bioinformatics Institute (EBI)

398    (https://www.ebi.ac.uk/Tools/msa/clustalo) and MUSCLE v3.8.31 [33] and examined using

399    Mesquite, a software for evolutionary biology (http://mesquiteproject.wikispaces.com/) [40].

400    Misaligned areas were manually adjusted after the MUSCLE alignment. Gaps and insertions in

401    the PMSA were analyzed to determine if the sequences in question were likely true indels or

402    likely to be artifacts of computer analysis of genome annotation.  BLAST searches were

403    performed of inserted runs of AAs that did not align with any other species in our PMSA, using

404    Protein BLAST, with default settings and query sequences of minimum length 30.  For a

405    "positive BLAST", the sequence results needed to show the presence of either homologs of the

406    query sequence in APC from other organisms or from known protein domains. For a "negative

407    BLAST", the only result was the sequence from the species used in the search query.  Exon

408    boundaries were identified using the NCBI Gene Database.  If an entire exon from one species

409    did not align with the other sequences and was deemed BLAST negative, that exon was

410    removed from the PMSA, using the rationale that it would be irrelevant to a variant found in

411    humans.

412    Phylogenetic trees were constructed from the curated APC alignment using a Maximum

413    Parsimony-based method implemented in PAUP* (Phylogenetic Analysis Using Parsimony

414    [*and Other Methods]), Version 4, Maximum Likelihood [34, 39], and Bayesian method as

415    implemented in MrBayes [35].

416    Nucleotide regions flanking prospective indels were analyzed using two splice site calculators:

417    (1) SpliceSiteFrame, (http://ibis.tau.ac.il/ssat/SpliceSiteFrame.htm), a splice site calculator from

418    Tel Aviv University, and (2) the online tool from the GENIE program [36]

419    (http://rulai.cshl.edu/new_alt_exon_db2/HTML/score.html), The maximum 3' score for a perfect

420    splice site would be 14.2, and the score for a perfect 5' splice score would be 12.6; these rarely

421    occur. Average scores for the 3' and 5' sites are 7.9 and 8.1 respectively.

422    *Substitution per site*

423    Absolute conservation of an amino acid in a PMSA can be determined with statistical

424    significance (P<0.05) if the PMSA contains at least three substitutions per site (subs/site, i.e.,

425    three times as many variants among all sequences as there are codons in the gene [6, 7]. In

426    order to determine if APC alignments contained three subs/site, we used the PHYLIP

427    (Phylogeny Inference Package) version 3.6a2 ProtPars program form the University of

428    Washington, Department of Genetics (http://evolution.genetics.washington.edu/phylip.html), with

429    the alignment converted to PHYLIP format. To convert the alignment from Clustal Omega

430    format to PHYLIP format and all other formats used during the analyses, the EMBOSS Seqret

431    from EBI (https://www.ebi.ac.uk/Tools/sfc/emboss_seqret/) and Mesquite Version 3.51 tools

432    were used (https://www.mesquiteproject.org/).

### *Predictions of Effects of APC Missense Substitutions*

434    In July 2013, 46 APC missense variants were collected from the LOVD database maintained by

435    the International Society for Gastrointestinal Hereditary Tumors (InSiGHT).  On May 30, 2018,

436    4891 variants observed by clinical genetic testing were collected from the ClinVar database

437    (http://www.ncbi.nlm.nih.gov/clinvar/).

438    **Computational Algorithms**:  The pathogenicity of each missense variant recorded in ClinVar

439    was predicted using the programs Align-GVGD, SIFT, PolyPhen2 MAPP, and REVEL.

440    *AlignGVGD* uses PMSAs and the biophysical properties of amino acid substitutions to calculate

441    the range of variation at each position.  Each variant is assigned a grade of C65 to C0

442    representing decreasing probability of deleterious, with C0 representing likely neutral AA

443    substitutions. [37] (http://agvgd.hci.utah.edu/about.php).

444    *SIFT (Sorting Intolerant From Tolerant)* creates position specific scoring matrices derived from

445    PMSAs. Each missense substitution predicted as "Tolerated' or "Affects Protein Function" [19].

446    (http://sift.bii.a-star.edu.sg/).

447    *PolyPhen2* combines its own pre-built sequence alignment with protein structural

448    characteristics, calculating a score used to classify each variant into three categories: benign,

449    possibly damaging and probably damaging. (http://genetics.bwh.harvard.edu/pph2/index.shtml)

450    [20]. We combined the categories of "possibly damaging" and "probably damaging".

451    *MAPP (Multivariate Analysis of Protein Polymorphisms)* also combines a PMSA with the

452    physiochemical characteristics of each AA position, predicting which AA should be deleterious

453    and which should be neutral at each position in the PMSA [38]

454    (http://www.ngrl.org.uk/Manchester/page/mapp-multivariate-analysis-protein-polymorphism).

455     ***REVEL** (Rare Exome Variant Ensemble Learner)* [16] is an ensemble method that uses

456     machine learning to combine the results of 13 individual predictors, using independent test sets

457     that did not overlap with sets used to train its component features. REVEL output classes were

458     designated as "Deleterious" for variants with a REVEL score ≥ 0.5 and "Neutral" with a REVEL

459     score < 0.5 [16].

460

461

462     **ACKNOWLEDGMENTS:**

466

467     **DISCLOSURES OF CONFLICT OF INTEREST:**

468     MER and TP are employees of Ambry Genetics, Inc.

469    REFERENCES:

470    1.  Hermel, D., McKinnon, W., Wood, M. and Greenblatt, M. (2016). Multi-gene panel testing for

471        hereditary cancer susceptibility in a rural Familial Cancer Program. *Familial Cancer*, 16(1),

472        pp.159-166.

473    2.  Yurgelun M, Allen B, Kaldate R, Bowles K, Judkins T, Kaushik P et al. Identification of a

474        Variety of Mutations in Cancer Predisposition Genes in Patients With Suspected Lynch

475        Syndrome. Gastroenterology. 2015;149(3):604-613.e20.

476    3.  Landrum M, Lee J, Benson M, Brown G, Chao C, Chitipiralla S et al. ClinVar: improving

477        access to variant interpretations and supporting evidence. Nucleic Acids Research.

478        2017;46(D1):D1062-D1067.

479    4.  Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J et al. Standards and guidelines

480        for the interpretation of sequence variants: a joint consensus recommendation of the

481        American College of Medical Genetics and Genomics and the Association for Molecular

482        Pathology. Genetics in Medicine. 2015;17(5):405-423.

483    5.  Tavtigian S, Greenblatt M, Lesueur F, Byrnes G. In silico analysis of missense substitutions

484        using sequence-alignment based methods. Human Mutation. 2008;29(11):1327-1336.

485    6.  Greenblatt M, Beaudet J, Gump J, Godin K, Trombley L, Koh J et al. Detailed computational

486        study of p53 and p16: using evolutionary sequence analysis and disease-associated

487        mutations to predict the functional consequences of allelic variants. Oncogene.

488        2003;22(8):1150-1163.

489    7.  Cooper G, Brudno M. Quantitative Estimates of Sequence Divergence for Comparative

490        Analyses of Mammalian Genomes. Genome Research. 2003;13(5):813-820.

491    8.  Thompson B, Greenblatt M, Vallee M, Herkert J, Tessereau C, Young E et al. Calibration of

492        Multiple In Silico Tools for Predicting Pathogenicity of Mismatch Repair Gene Missense

493        Substitutions. Human Mutation. 2012;34(1):255-265.

494   9.  Tavtigian S, Samollow P, de Silva D, Thomas A. An Analysis of Unclassified Missense

495       Substitutions in Human BRCA1. Familial Cancer. 2006;5(1):77-88.

496   10. Goldgar D, Easton D, Byrnes G, Spurdle A, Iversen E, Greenblatt M. Genetic evidence and

497       integration of various data sources for classifying uncertain variants into a single model.

498       Human Mutation. 2008;29(11):1265-1272.

499   11. Abkevich V, Zharkikh A, Deffenbaugh AM, Frank D, Chen Y, Shattuck D, Skolnick MH, Gutin

500       A, Tavtigian SV. Analysis of missense variation in human BRCA1 in the context of

501       interspecific sequence variation. Journal of Medical Genetics. 2004;41(7):492-507.

502   12. Miller P, Duraisamy S, Newell J, Chan P, Tie M, Rogers A et al. Classifying variants of

503       CDKN2A using computational and laboratory studies. Human Mutation. 2011;32(8):900-911.

504   13. Rishishwar L, Varghese N, Tyagi E, Harvey S, Jordan I, McCarty N. Relating the Disease

505       Mutation Spectrum to the Evolution of the Cystic Fibrosis Transmembrane Conductance

506       Regulator (CFTR). PLoS ONE. 2012;7(8):e42336.

507   14. Scott R, Crooks R, Rose L, Attia J, Thakkinstian A, Thomas L et al. Germline Missense

508       Changes in the APC Gene and Their Relationship to Disease. Hereditary Cancer in Clinical

509       Practice. 2004;2(2):81.

510   15. Grandval P, Blayau M, Buisine M, Coulet F, Maugard C, Pinson S et al. The UMD-APC

511       Database, a Model of Nation-Wide Knowledge Base: Update with Data from 3,581

512       Variations. Human Mutation. 2014;35(5):532-536.

513   16. Ioannidis N, Rothstein J, Pejaver V, Middha S, McDonnell S, Baheti S et al. REVEL: An

514       Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants. The

515       American Journal of Human Genetics. 2016;99(4):877-885.

516   17. Tavtigian, S., Deffenbaugh, Yin, Judkins, Scholl, Samollow, De Silva, Zharkikh, Thomas.

517       Comprehensive statistical study of 452 BRCA1 missense substitutions with classification of

518       eight recurrent substitutions as neutral. Journal of Medical Genetics. 2005;43(4):295-305.

519  18. Tavtigian S, Greenblatt M, Harrison S, Nussbaum R, Prabhu S, Boucher K et al. Modeling

520      the ACMG/AMP variant classification guidelines as a Bayesian classification framework.

521      Genetics in Medicine. 2018;20(9):1054-1060.

522  19. Ng PC, Henikoff S. SIFT: predicting amino acid changes that affect protein function. Nucleic

523      Acids Research. 2003;31(13):3812-3814.

524  20. Adzhubei I, Jordan D, Sunyaev S. Predicting Functional Effect of Human Missense

525      Mutations Using PolyPhen-2. Current Protocols in Human Genetics. 2013;76(1):7.20.1-

526      7.20.41.

527  21. Aoki K, Taketo M. Adenomatous polyposis coli (APC): a multi-functional tumor suppressor

528      gene. Journal of Cell Science. 2007;120(19):3327-3335.

529  22. Huber A, Weis W. The Structure of the β-Catenin/E-Cadherin Complex and the Molecular

530      Basis of Diverse Ligand Recognition by β-Catenin. Cell. 2001;105(3):391-402.

531  23. Menéndez M, González S, Obrador–Hevia A, Domínguez A, Pujol M, Valls J et al.

532      Functional Characterization of the Novel APC N1026S Variant Associated With Attenuated

533      Familial Adenomatous Polyposis. Gastroenterology. 2008;134(1):56-64.

534  24. Grimm D, Azencott C, Aicheler F, Gieraths U, MacArthur D, Samocha K et al. The

535      Evaluation of Tools Used to Predict the Impact of Missense Variants Is Hindered by Two

536      Types of Circularity. Human Mutation. 2015;36(5):513-523.

537  25. Drost M, Tiersma Y, Thompson B, Frederiksen J, Keijzers G, Glubb D et al. A functional

538      assay–based procedure to classify mismatch repair gene variants in Lynch syndrome.

539      Genetics in Medicine. 2018;.

540  26. Adebali O, Reznik A, Ory D, Zhulin I. Establishing the precise evolutionary history of a gene

541      improves prediction of disease-causing missense mutations. Genetics in Medicine.

542      2016;18(10):1029-1036.

543  27. Ortuño F, Valenzuela O, Rojas F, Pomares H, Florido J, Urquiza J et al. Optimizing multiple

544      sequence alignments using a genetic algorithm based on three objectives: structural

545    information, non-gaps percentage and totally conserved columns. Bioinformatics.

546    2013;29(17):2112-2121.

547    28. Lee K, Krempely K, Roberts M, Anderson M, Carneiro F, Chao E et al. Specifications of the

548    ACMG/AMP variant curation guidelines for the analysis of germline CDH1 sequence

549    variants. Human Mutation. 2018;39(11):1553-1568.

550    29. Goldgar D, Easton D, Deffenbaugh A, Monteiro A, Tavtigian S, Couch F. Integrated

551    Evaluation of DNA Sequence Variants of Unknown Clinical Significance: Application to

552    BRCA1 and BRCA2. The American Journal of Human Genetics. 2004;75(4):535-544.

553    30. Easton D, Deffenbaugh A, Pruss D, Frye C, Wenstrup R, Allen-Brady K et al. A Systematic

554    Genetic Assessment of 1,433 Sequence Variants of Unknown Clinical Significance in the

555    BRCA1 and BRCA2 Breast Cancer–Predisposition Genes. The American Journal of Human

556    Genetics. 2007;81(5):873-883.

557    31. Hansford S, Kaurah P, Li-Chang H, Woo M, Senz J, Pinheiro H et al. Hereditary Diffuse

558    Gastric Cancer Syndrome. JAMA Oncology. 2015;1(1):23.

559    32. Dommering C, Mol B, Moll A, Burton M, Cloos J, Dorsman J et al. RB1mutation spectrum in

560    a comprehensive nationwide cohort of retinoblastoma patients. Journal of Medical Genetics.

561    2014;51(6):366-374.

562    33. Edgar R. MUSCLE: multiple sequence alignment with high accuracy and high throughput.

563    Nucleic Acids Research. 2004;32(5):1792-1797.

564    34. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large

565    phylogenies. Bioinformatics. 2014;30(9):1312-1313.

566    35. Ronquist F, Huelsenbeck J. MrBayes 3: Bayesian phylogenetic inference under mixed

567    models. Bioinformatics. 2003;19(12):1572-1574.

568    36. Reese M, Eeckman F, Kulp D, Haussler D. Improved Splice Site Detection in Genie. Journal

569    of Computational Biology. 1997;4(3):311-323.

570   37. Tavtigian SV, Byrnes GB, Goldgar DE, Thomas A. Classification of rare missense

571        substitutions, using risk surfaces, with genetic-and molecular-epidemiology applications

572        Hum Mutat. 2008 Nov;29(11):1342-54.

573   38. Stone EA, Sidow A. Physicochemical constraint violation by missense substitutions

574        mediates impairment of protein function and disease severity. Genome Research.

575        2005;15(7):978-986.

576   39. Swofford, D. L. 2002. PAUP*. Phylogenetic Analysis Using Parsimony (*and Other

577        Methods). Version 4. Sinauer Associates, Sunderland, Massachusetts.

578   40. Maddison, W. P. and D.R. Maddison. 2018. Mesquite: a modular system for evolutionary

579        analysis. Version 3.51  http://www.mesquiteproject.org

580    **FIGURES**

581

```
Sea Squirt   AESIIHSLREIMELVNYKVQRLAINELGGLFCVAEILILHCSSKHDEEAQEETGSRLRQY   419
Sea Urchin   --NPGPAMAALMKLSFDEEHRSAICHLGGLHAIAELLQVDYEVHGS--SNDQYTVTLRRY   442
Zebrafish    --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIGELLQVDCEIYGL--TNDHYSVTLRRY   494
Frog         --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--INDHYSVTLRRY   502
Chicken      --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--TNDHYSVTLRRY   500
Opossum      --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--TSDHYSVTLRRY   500
Mouse        --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--TNDHYSVTLRRY   498
Cow          --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--TNDHYSITLRRY   500
Monkey       --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--TNDHYSITLRRY   500
Human        --QICPAVCVLMKLSFDEEHRHAMNELGGLQAIAELLQVDCEMYGL--TNDHYSITLRRY   500
             .    ::   :*:*    : :* *: .**** .:.*:* :. .       .:.    **:*

Sea Squirt   SGRILTNLTYADNLNKVLLMNMRGLLETVRDQLQHESEEIQQAMASILRNLSWQADKEGR   479
Sea Urchin   AGMALTNLTFGDVTNKALLCSMKGCMKALVALLSAESEDLRQVAASVLRNLSWRADMASK   502
Zebrafish    AGMALTNLTFGDVANKATLCSMKGCMRAMVAQLKSESEDLQQVIASVLRNLSWRADVNSK   554
Frog         AGMALTNLTFGDVANKATLCSMKSCMRALVAQLKSESEDLQQVIASVLRNLSWRADVNSK   562
Chicken      AGMALTNLTFGDVANKATLCSMKGCMRALVAQLKSESEDLQQVIASVLRNLSWRADVNSK   560
Opossum      AGMALTNLTFGDVANKATLCSMKGCMRALVAQLKSESEDLEQVIASVLRNLSWRADVNSK   560
Mouse        AGMALTNLTFGDVANKATLCSMKGCMRALVAQLKSESEDLQQVIASVLRNLSWRADVNSK   558
Cow          AGMALTNLTFGDVANKATLCSMKGCMRALVAQLQSESEDLQQVIASVLRNLSWRADVNSK   560
Monkey       AGMALTNLTFGDVANKATLCSMKGCMRALVAQLKSESEDLQQVIASVLRNLSWRADVNSK   560
Human        AGMALTNLTFGDVANKATLCSMKGCMRALVAQLKSESEDLQQVIASVLRNLSWRADVNSK   560
             :*   *****:.*   **. * .*:. :.::   *. ***::.*. **:******:**   .:
```

582

583    **Figure 1. Excerpt of the curated APC alignment generated from the MSA program Clustal**

584    **Omega.**

585    Exon boundaries are labeled in red with a black background. The red highlighted region in the

586    human sequence corresponds to a portion of an Armadillo Repeat domain.
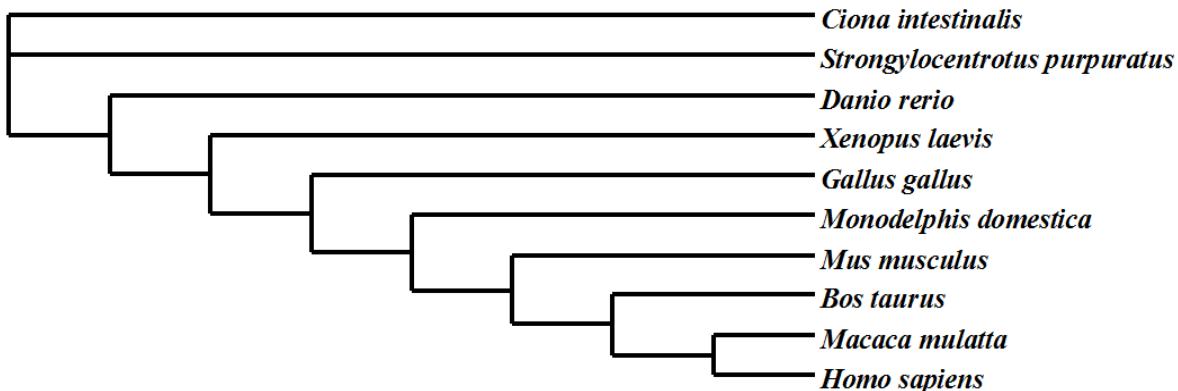
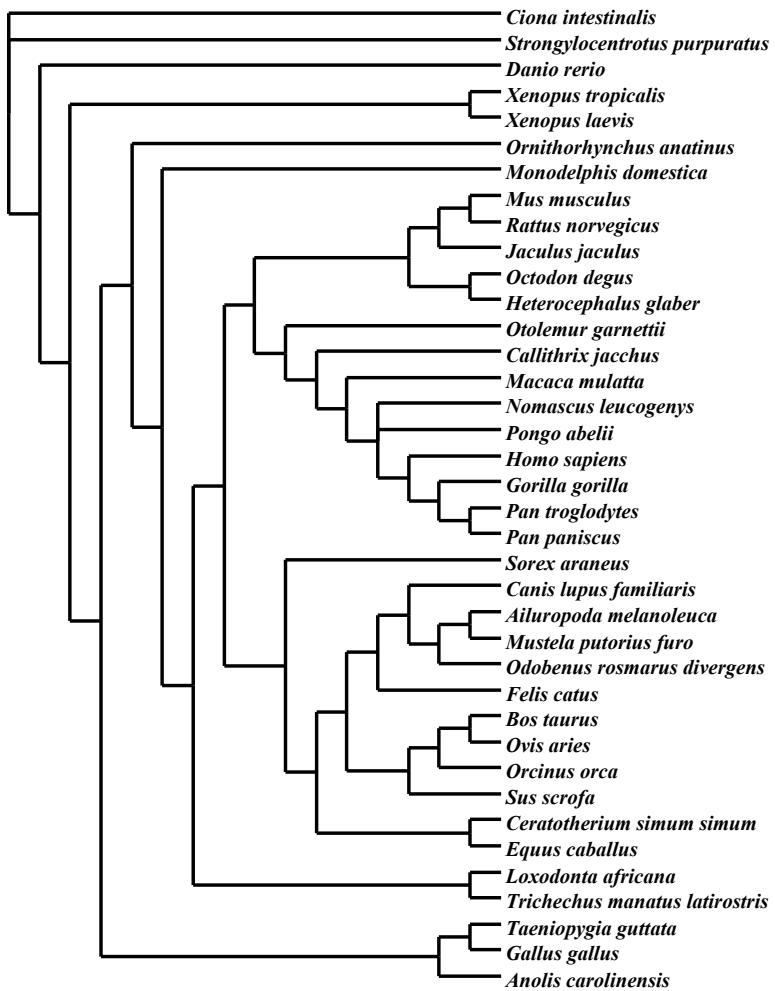587

588

589

590

591

592

593

594

595

596

597

**Figure 2A. Ten species phylogenetic consensus tree for the APC protein constructed using the computational phylogenetics program PAUP\* (Phylogenetic Analysis Using Parsimony \*and other methods).**

601

602

**Figure 2B. Thirty-eight species phylogenetic consensus tree for the APC protein**

**constructed using the computational phylogenetics program PAUP\* (Phylogenetic**

**Analysis Using Parsimony \*and other methods).**

606

607

608

609

610

611

612

613     **(PDF in Separate File)**

614     **Supplementary Figure 1. Curated 10-Species APC alignment.** PMSA was generated from

615     the program Clustal Omega. Exon boundaries are labeled in red with a black background. The

616     domains are highlighted throughout the alignment. Grey is oligomerization domain, red is

617     Armadillo repeats, yellow is Beta Catenin Repeats, green is a sequence with homology to the

618     herpes virus (PHA03307), turquoise is the Basic domain, and purple is the EB1 and HDLG

619     binding site.

620

621

622     **(PDF in Separate File)**

623     **Supplementary Figure 2. 38-Species APC alignment.** PMSA was generated from the

624     program Clustal Omega.  No annotation is added.

625

626