# Independent evolution of tetraloop in enterovirus oriL replicative element and its putative binding partners in protein 3C

**Maria A Prostova** [Corresp., 1] , **Andrey A Deviatkin** [1] , **Irina O Tcelykh** [1, 2] , **Alexander N Lukashev** [1, 3] , **Anatoly P Gmyl** [1, 2, 3]

[1] Chumakov Institute of Poliomyelitis and Viral Encephalitides, Moscow, Russia

[2] Lomonosov Moscow State University, Moscow, Russia

[3] Sechenov First Moscow State Medical University, Moscow, Russia

Corresponding Author: Maria A Prostova
Email address: prostova_ma@chumakovs.su

**Background.** Enteroviruses are small non-enveloped viruses with (+) ssRNA genome with one open reading frame. Enterovirus protein 3C (or 3CD for some species) binds replicative element oriL to initiate replication. The replication of enteroviruses features low fidelity that allows virus to adapt to changing environment on the one hand, and requires additional mechanisms to maintain the genome stability on the other. Structural disturbances in the apical region of oriL domain d can be compensated by amino acid substitutions in positions 154 or 156 of 3C (amino acid numeration corresponds to poliovirus 3C), thus suggesting co-evolution of these interacting sequences in nature. The aim of this work was to understand co-evolution patterns of two interacting replication machinery elements in enteroviruses, the apical region of oriL domain d and its putative binding partners in the 3C protein.

**Methods.** To evaluate the variability of the domain d loop sequence we retrieved all available full enterovirus sequences (>6400 nucleotides) that were present in the NCBI database on February 2017 and analyzed variety and abundance of sequences in domain d of replicative element oriL and in the protein 3C.

**Results.** A total of 2842 full genome sequences were analyzed. Majority of domain d apical loops were tetraloops, which belonged to consensus YNHG (Y=U/C, N=any nucleotide, H=A/C/U). Putative RNA-binding tripeptide 154-156 (Enterovirus C 3C protein numeration) was less diverse than the apical domain d loop region and, in contrast to it, was species-specific.

**Discussion.** Despite RNA-binding tripeptide is suggested to interact with apical region of domain d, they evolve independently in nature. Together, our data indicates plastic evolution of both interplayers of 3C-oriL recognition.

1 Independent evolution of tetraloop in enterovirus oriL

2 replicative element and its putative binding partners in

3 virus protein 3C.

4 Maria A. Prostova[1], Andrey A. Deviatkin[1], Irina O. Tcelykh[1,2], Alexander N.

5 Lukashev[1,3], Anatoly P. Gmyl[1,2,3]

6 1 - Chumakov Institute of Poliomyelitis and Viral Encephalitides, Moscow, Russian

7 Federation

8 2 - Lomonosov Moscow State University, Moscow, Russia

9 3 - Sechenov First Moscow State Medical University, Moscow, Russia.

10

11 Corresponding author:

12 Maria A. Prostova[1]

13 Email address:

14 prostovna@gmail.com

15

# Abstract

16

17      **Background.** Enteroviruses are small non-enveloped viruses with (+) ssRNA genome

18   with one open reading frame. Enterovirus protein 3C (or 3CD for some species) binds replicative

19   element oriL to initiate replication. The replication of enteroviruses features low fidelity that

20   allows virus to adapt to changing environment on the one hand, and requires additional

21   mechanisms to maintain the genome stability on the other. Structural disturbances in the apical

22   region of oriL domain d can be compensated by amino acid substitutions in positions 154 or 156

23   of 3C (amino acid numeration corresponds to poliovirus 3C), thus suggesting co-evolution of

24   these interacting sequences in nature. The aim of this work was to understand co-evolution

25   patterns of two interacting replication machinery elements in enteroviruses, the apical region of

26   oriL domain d and its putative binding partners in the 3C protein.

27

28      **Methods.** To evaluate the variability of the domain d loop sequence we retrieved all

29   available full enterovirus sequences (>6400 nucleotides) that were present in the NCBI database

30   on February 2017 and analyzed variety and abundance of sequences in domain d of replicative

31   element oriL and in the protein 3C.

32

33      **Results.** A total of 2842 full genome sequences were analyzed. Majority of domain d

34   apical loops were tetraloops, which belonged to consensus YNHG (Y=U/C, N=any nucleotide,

35    H=A/C/U). Putative RNA-binding tripeptide 154-156 (Enterovirus C 3C protein numeration)

36    was less diverse than the apical domain d loop region and, in contrast to it, was species-specific.

37

38    **Discussion.** Despite RNA-binding tripeptide is suggested to interact with apical region of

39    domain d, they evolve independently in nature. Together, our data indicates plastic evolution of

40    both interplayers of 3C-oriL recognition.

41

42                                    **Introduction**

43            Enteroviruses are small non-enveloped viruses with plus strand genome about 7500 nt,

44    which contains one open reading frame that encodes structural (capsid) and non-structural

45    proteins, 5' and 3' NTRs (not translated regions) and polyA on the 3' end (Palmenberg,

46    Neubauer and Skern, 2010). (Figure 1). Most non-structural enterovirus proteins are

47    polyfunctional. Protease 3CD is a precursor of polymerase 3D and plays the key role in the

48    initiation of replication (Harris *et al.*, 1992; Gamarnik and Andino, 1998; Thompson and

49    Peersen, 2004). After translation by host cell ribosomal machinery, the genome is utilized for the

50    synthesis of the (-) strand RNA, which, in turn, serves as a matrix for synthesis of multiple

51    daughter (+) strands. Non-translated regions of genome and a coding sequence within the

52    genomic region encoding the viral helicase 2C contain replicative elements, which interact with

53    viral and host proteins. These RNA-protein complexes regulate initiation and further steps of the

54    replication. For poliovirus, the most clinically relevant member of the Enterovirus genus, there

55    are at least three known RNA-protein complexes, which are formed with replicative elements

56    oriL, oriR and oriI during replication (Figure 1).

57        Complex of oriL with viral protein 3CD and the host protein PCBP2 is crucial for the

58    transcription initiation (Goodfellow *et al.*, 2000; Vogt and Andino, 2010; Chase, Daijogo and

59    Semler, 2014). Element oriL has a cloverleaf-like secondary structure with four domains termed

60    "a" (stem of the cloverleaf), "b", "c" and "d" (leafs of the cloverleaf) (Trono, Andino and

61    Baltimore, 1988; Andino, Rieckhof and Baltimore, 1990) (Figure 1). Previously, it was

62    demonstrated *in vitro* that 3CD (or 3C) of poliovirus, coxsackievirus B3 and bovine enterovirus

63    1 interact with the apical loop and the flanking base pairs of hairpin d in the oriL element

64    (Andino, Rieckhof and Baltimore, 1990; Du *et al.*, 2003; Ihle *et al.*, 2005) (Figure 1).

65        The apical loops of domain d in genomes belonging to several viruses of *Enterovirus*

66    genera were shown by NMR experiments to be tetraloops with a specific spatial structure, which

67    belongs to the UNCG structural class of stable tetraloops (Du *et al.*, 2003, 2004; Ihle *et al.*, 2005;

68    Melchers *et al.*, 2006). There are several known structural classes of tetraloops and three of

69    them, named according to consensus sequences, contain tetraloops of extremal stability: UNCG

70    (where N=any nucleotide), GNRA (where R=A/G) and gCUUGc (Uhlenbeck, 1990; Cheong and

71    Cheong, 2010). Tetraloops of UNCG and GNRA classes are the most widely represented (Woese

72    *et al.*, 1990; Cheong and Cheong, 2010; Bottaro and Lindorff-Larsen, 2017).  Previously it was

73    shown that only tetraloops of UNCG structural class, but not tetraloops of GNRA or gCUUGc

74    structural classes can support effective replication of poliovirus genome (Prostova *et al.*, 2015).

75    Moreover, the exact sequence of the apical region of poliovirus domain d was of less importance

76    for effective 3CD-oriL recognition than its spatial structure (Rieder *et al.*, 2003; Prostova *et al.*,

77    2015). At the same time, structural disturbance in the apical region of oriL domain d of

78    poliovirus could be compensated by amino acid substitutions in the tripeptide 154-156 of the 3C

79    protein (here and after amino acid numeration corresponds to poliovirus 3C protein) (Andino *et*

80    *al.*, 1990; Prostova *et al.*, 2015). In addition to triplet 154-156, the conservative motif $_{82}$KFRDI$_{86}$

81    of the 3C protein also takes part in the oriL recognition (Andino *et al.*, 1990, 1993; Hämmerle,

82    Molla and Wimmer, 1992; Shih, Chen and Wu, 2004).

83        The replication of enterovirus is a low fidelity process, generating on average one

84    mutation per genome (Sanjuán *et al.*, 2010; Acevedo, Brodsky and Andino, 2013). The high

85    probability of mutation allows virus to adapt to constantly changing environment on the one

86    hand, but requires additional mechanisms to maintain the genome stability on the other (Wagner

87    and Stadler, 1999; Lauring, Frydman and Andino, 2013). The aim of the present study was to

88    understand co-evolution patterns of two interacting replication machinery elements in

89    enteroviruses, the apical domain d of oriL and the 3C protein.

90

91                                    **Materials and methods**

92                        **Formation and filtration of sets of full genomes**

93            Genomes of species of *Enterovirus* genus with length 8000>n>6800 were extracted from

94     the NCBI database. For every species, a multiple sequence alignment was done with Clustal

95     (Larkin *et al.*, 2007). Sequences that contained more than 50 "N" characters in succession and

96     sequences that were annotated as "Modified_Microbial_Nucleic_Acid", were deleted from

97     alignments. In order to reduce the bias towards particular loop sequences present in a large set of

98     closely related genomes, which, for example, belonged to one outbreak, all sequences that

99     differed from any other sequence in the dataset by less than 1% of the nucleotide sequence were

100    omitted.

101                    **Analysis of tetraloop and amino acid variety in the sets of genomes**

102            For analysis of domain d sequence variety, the multiple sequence alignments were used.

103    The relevant region of multiple sequence alignment and respective names of sequences were

104    analyzed in Microsoft Excel. To analyze correlation of domain d loop and tripeptide of 3C

105    sequences the same alignments were translated in the protein 3C coding region. The resulting

106    amino acid sequences that corresponded to tripeptides 154-156 (poliovirus 3C numerations) were

107    analyzed using Microsoft Excel. Amino acid frequency plot was formed with WebLogo server

108    using the set of filtered genomes for every species (Crooks *et al.*, 2004). To do this the multiple

109    sequence alignment of filtered genomes of every species was translated in the region that codes

110    protein 3C, and positions 71-89 and 147-160 were saved in separate MAS files, which were then

111    used to produce Logo.

## Domain d secondary structure

113        Domain d secondary structure was folded with Vienna RNA Websuit server with

114    subsequent manual editing (Gruber *et al.*, 2008; Lorenz *et al.*, 2011).

# Results

## Sample characteristics

117        To evaluate the variability of the domain d loop sequence we retrieved all available

118    complete genome (8000>n>6800 nucleotides) enterovirus (EV) sequences that were present in

119    the NCBI database on February 2017. A total of 1173, 414, 773, 462, 12, 13, 15, 3, 7, 202, 76

120    and 51 sequences were analyzed for species *Enterovirus A, Enterovirus B, Enterovirus C,*

121    *Enterovirus D, Enterovirus E, Enterovirus F, Enterovirus G, Enterovirus H, Enterovirus J,*

122    *Rhinovirus A, Rhinovirus B* and *Rhinovirus C*, respectively. As expected, genomes of

123    epidemiologically significant viruses were the most represented in database. For example, 66%

124    of *Enterovirus A* species genomes belonged to EV71 type, the causative agent of hand, foot and

125    mouth disease (Solomon *et al.*, 2010); most *Enterovirus C* species sequences (78%) belonged to

126     poliovirus; most *Enterovirus D* species sequences (98.7%) represented EV68, an etiological

127     agent of severe respiratory illness (Oermann *et al.*, 2015). The number of genome sequences of

128     each species that contained the oriL region is shown in Table 1.

129        Sequences of apical regions in oriL domain d and the amino acids involved in RNA

130     recognition in 3C protein were analyzed. In genomes of *Enterovirus E* and *Enterovirus F* species

131     that have two oriLs (Pilipenko, Blinov and Agol, 1990; Zell *et al.*, 1999) sequences of both oriLs

132     were analyzed (Table 1). To reduce the bias towards particular loop sequences present in a large

133     set of closely related genomes, which, for example, belonged to one outbreak, all sequences that

134     differed from any other sequence in the dataset by less than 1% of the nucleotide sequence were

135     omitted. After curation, the sizes of the largest data sets decreased dramatically, but the number

136     of unique loop sequences in every set did not change significantly (Table 1). Unique tetraloop

137     variants were lost for *Enterovirus A* (tetraloop UGUG), Enterovirus C (tetraloops CCCG, CAUG

138     and UGUG) and *Enterovirus D* (tetraloop UUGG). This indicates that even among closely

139     related genomes tetraloop sequence can vary. Indeed, in several outbreaks caused by EV71 or

140     PV1 closely related genomes contained different apical domain d sequences (not shown). It

141     should be noted that filtration of the dataset using a 95% sequence identity threshold resulted in a

142     dramatic loss of unique tetraloop variants (data not shown).

143                  **Variability of oriL domain d apical loop sequence**

144        The secondary structure of domain d was conservative in all species, except *Enterovirus*

145     *G*, which has elongated domain d (Figure 2) (Krumbholz *et al.*, 2002).

146        Variety and occurrence of various loops in apical region of domain d in all species of

147     Enterovirus genera were analyzed in filtered sets of full genome sequences. Most of domain d

148     apical loops were tetraloops (i.e. consisted of four nucleotides) (Table 2). However, triloops (3-

149     nucleotide loop) could be found in genomes of *Enterovirus C* and *Rhinovirus A* and *B* species,

150     whereas pentaloops (5-nucleotide loop) were detected in genomes of *Enterovirus E* species

151     (Table 2).

152        The most common loop sequences belonged to consensuses YNMG (Y=C/U, N=any,

153     M=A/C; tetraloops with UNCG class spatial structure belong to this consensus) and YNUG

154     (tetraloops with UNCG class and gCUUGc class spatial structures belong to this consensus)

155     (Table 2). Consensus YNMG and consensus YNUG together corresponded to 24 unique

156     sequence variants. Interestingly, in our dataset of 2842 full genomes, four tetraloops out of 24

157     possible variants have never been found in the domain d apical region: UUAG, UCAG, CUAG

158     and CCAG (Table 2). Thus, dinucleotides UA and CA are likely to be avoided at positions 2-3 of

159     tetraloop in enterovirus genomes.

160        In enterovirus A species 17 out of 24 possible unique tetraloop sequences were identified

161     (Table 2, Table S2). Twelve unique loops of *Enterovirus A* belonged to consensus YNMG, while

162     5 others belonged to consensus YNUG. Frequency of particular tetraloop sequences varied

163     significantly (Table 2, Table S2). One tetraloop (UGUG) was lost upon filtration and manually

164     added to the final data set (Table 2). Interestingly, the diversity of tetraloops among EV71

165     serotype was similar to the diversity of tetraloops in the whole *Enterovirus A* species (Table S1).

166          In *EV-B* genomes 18 unique tetraloops out of 24 possible were found. Twelve of these

167     tetraloops belonged to consensus YNMG and six to consensus YNUG (Table 2, Table S3). The

168     most abundant tetraloops were CACG (98 genomes), UACG (51 genomes) and UGCG (31

169     genomes), that also were present among the most abundant tetraloops of species Enterovirus A.

170          In genomes of *Enterovirus C* species 9 unique tetraloops belonged to YNMG consensus

171     and 4 to YNUG consensus. Three unique tetraloops were lost upon filtration and added to final

172     data set (CCCG, UGUG, CAUG) (Table 2, Table S4). Two genomes annotated at the NCBI data

173     base as Human coxsackievirus A21 strain Coe (accession number D00538) and Human

174     coxsackievirus A21 strain BAN00-10467 (accession number EF015031) contained triloops CAG

175     and CCG, respectively. The most abundant tetraloops in *EV-C* species were UACG (106),

176     CACG (101), UGCG (43) (Table 2, Table S4), which corresponds to the Sabin vaccine strains of

177     poliovirus serotypes 2, 3 and 1, respectively. To evaluate bias caused by redundant number of

178     vaccine strain sequences in the data set, we subtracted genomes of vaccine/vaccine derived

179     poliovirus strains from the analyzed set. As a result, tetraloops UACG, CACG and UGCG were

180     still the most frequent variants (Table S1).

181          Just 57 *EV-D* genomes out of 419 were left after 1%-identity filtration. Fifty genomes

182     belonged to Human enterovirus 68, the etiological agent of respiratory illness. All genomes of

183    this type contained loop UUCG in the domain d apical region. Other tetraloops were UUUG (1),

184    CUCG (2), CCCG (1), CUUG (2) and CACG (1) (Table 2, Table S5). One tetraloop (UUGG)

185    was lost upon filtration and manually added to the final data set.

186         Species *Enterovirus E* and *F* have 2 oriLs in the 5'NTR, with, in general, similar

187    sequences in apical region of domain d (Pilipenko, Blinov and Agol, 1990; Zell *et al.*, 1999)

188    (Table S6). Due to this, we united sequences from the first and the second oriL of these viruses

189    in heat map (Table 2). Domain d loops in 10 genomes of *EV-F* were tetraloops, while in 10 *EV-E*

190    genomes there were both tetraloops (first oriL) and pentaloops (the first and the second oriL)

191    (Table 2, Table S6). There were four diverse tetraloop sequences in oriLs of Enterovirus E and F

192    with no obvious preference between these species. These sequences were GCUA, GUUA,

193    GCCA, AUUA (Table 2, Table S6). Tetraloop AUUA was found once in the first oriL domain d

194    of EV-F (strain PS87/Belfast, accession number DQ092794) (Table 2, Table S6). There were six

195    diverse pentaloop sequences in domain d of Enterovirus E genomes – GCUUA, GUUUA,

196    GCCUA, GCGUA, GAUUA, GUCUA (Table 2, Table S6).

197         All domain d loops in genomes of Enterovirus G, H and J species were tetraloops, and all

198    except one tetraloop variant belonged to consensus YNMG (Table 2, Table S7). One Enterovirus

199    G representative had GUUA tetraloop sequence (strain LP 54, accession number AF363455),

200    similar to loops of Enterovirus E and F species (Table 2).  This genome had only one oriL with

201    the same domain d length as of Enterovirus G genomes (Krumbholz *et al.*, 2002).

202     All except one (isolate V38_URT-6.3m, accession number JF285329) full genomes of

203     Rhinoviruses A species and all full genomes of *Rhinovirus C* species had tetraloops in apical

204     regions of domain d (Table 2). Tetraloops of these viruses in almost all cases belonged to

205     consensus YNMG, with one exception in *Rhinovirus C* (tetraloop CUUC, isolate JAL-1,

206     accession number JX291115) (Table 2, Table S8).  All loops in the *Rhinovirus B* domain d apical

207     region were triloops (Table 2, Table S8).

208     Thus, the secondary structure of domain d was very similar among species of the genus

209     *Enterovirus,* with the exception of *Enterovirus G* species (Figure 2). Apical region of domain d

210     has high diversity of sequences, however in species *Enterovirus A, B, C, D, G, H, J* and

211     *Rhinovirus A* and *C* it mostly can correspond to the same consensus YNHG (Y=C/U, N=any,

212     H=A/C/U).

213                         **Variety of RNA-recognition tripeptide of 3C**

214     Two motifs of protein 3C are involved in RNA recognition and interact with oriL: the

215     conservative motif KFRDI (positions 82-86 of poliovirus 3C) and putative RNA-binding

216     tripeptide (positions 154-156 in poliovirus 3C) (Andino *et al.*, 1990, 1993; Hämmerle, Hellen

217     and Wimmer, 1991; Shih, Chen and Wu, 2004). Substitutions in putative RNA-binding tripeptide

218     are known to compensate the disturbance in apical region of domain d, and therefore the RNA-

219     binding tripeptide is a putative candidate to co-evolve with domain d loop (Andino *et al.*, 1990;

220     Prostova *et al.*, 2015). There are other amino acids that are shown to affect oriL-3CD interaction,

221    but tripeptide 154-156 (*Enterovirus C* 3C protein numbering here and below) is the only one that

222    was proven to compensate structural disturbance in domain d apical region (Andino *et al.*, 1990,

223    1993). To evaluate the possible co-evolution between domain d tetraloop and its putative

224    interaction partners in the protein 3C, relevant sequences in the filtered full genome data sets

225    were analyzed.

226         Motif $_{82}$KFRDI$_{86}$ was conserved in all species, as well as amino acids Glu 71 and Cys

227    147 of the protease catalytic triade (Figure 3). Second position of the putative RNA-binding

228    tripeptide (position 155) was invariantly Gly.

229         No mutual dependence between loop sequences and tripeptide sequences was found

230    within enterovirus genomes of the same species. For example, *Enterovirus A* genomes contained

231    17 unique variants of tetraloop sequence, whereas the predominant fraction of 3C sequences (548

232    out of 564) contained the conservative tripeptide VGK at positions 154-156 (Figure 3, Table S9).

233    Noteworthy, this tripeptide was found not only in genomes of EV71 serotype, although genomes

234    of this serotype prevailed in the dataset. Other *EV-A* genomes contained tripeptides VGR (7 out

235    of 564), TGK (4 out of 564), IGK (3 out of 564), VGE (1 out of 564) and SRK (1 out of 564)

236    (Figure 3, Table S9). Genomes with tripeptides other than VGK did not contain any peculiarities

237    of domain d loop sequence (Table S9). This observation confirms that the specific loop sequence

238    is likely not the main subject for recognition by the RNA-binding tripeptide. Similarly, all or

239    almost all genomes of *Enterovirus B* (242 out of 244), *Enterovirus C* (272 out of 274),

240     *Enterovirus D* (all), *Enterovirus G* (7 out of 8), *Enterovirus H* (a total of 2 genomes – one with

241     TGK, one with TGR), *Enterovirus J* (all) and *Rhinovirus B* (36 out of 37) species contained

242     tripeptide TGK at position 154-156 of 3C protein (Figure 3, Table S7, S9, S10). Alternative

243     tripeptides were TGR in two genomes of *Enterovirus B* and one genome of *Enterovirus H*; IGK

244     in one genome of *Enterovirus C* species and in one genome of *Rhinovirus B* species; PGK in one

245     genome of *Enterovirus C* species; MGK in one genome of *Enterovirus G* species (Table S7, S9,

246     S10).

247          Genomes of species *Enterovirus E* and *F* contained two oriLs with tetraloops in domain d

248     mostly of consensus GYYA or pentaloops of consensus GHBUA, where H = A/C/U, and B =

249     U/C/G. All genomes contained tripeptide MGK at positions 154-156 of the protein 3C (Figure 3,

250     Table S7). Interestingly, a similar loop-tripeptide pair was found in one genome of *Enterovirus G*

251     species (strain LP 54, accession number AF363455). It contained tetraloop GUUA in the domain

252     d of its single oriL and tripeptide MGK in 3C. Unlike this unique genome, other genomes of

253     *Enterovirus G* species contained tetraloops of YNMG consensus and tripeptide TGK in the

254     protein 3C.

255          Rhinovirus genomes contained tetraloops mostly of consensus YNMG (species

256     *Rhinovirus A* and *C*) or triloops (*Rhinovirus A* and *B*) (Table 2). *Rhinovirus A* genomes with

257     tetraloops in domain d contained tripeptides in 3C with positively charged amino acid before the

258     tripeptide, but not at the last position of it, as in genomes of *Enterovirus A-C* species (Figure 3,

259    Table S11, S12). Sequence of tripeptides did not depend on tetraloop sequence and was, in

260    descending order, IGQ (the most abundant, 65 genomes out of 119), IGL (20 genomes out of

261    119), IGS, VGS, IGN, VGQ, IGV, VGH (Table S11). In the case of *Rhinovirus A* genome with

262    triloop UCU in domain d (isolate V38_URT-6.3m, accession number JF285329), protein 3C

263    contained tripeptide TGK without positive charged amino acid before it (Table S11). All

264    genomes of *Rhinovirus B* species contained triloops in the domain d and all, except one with

265    IGK, contained tripeptide TGK in 3C. Genomes of Rhinovirus C contained tetraloops mostly of

266    consensus YHCG (H=all but G) and tripeptides in 3C without a positively charged amino acid at

267    the last position (TGN, VGN, TGH) or outside of the tripeptide (Table S12). One genome

268    contained tetraloop CUUC paired with most abundant tripeptide TGN (23 out of 37

269    genomes)(Table S12).

270        Thus, dependence between apical domain d sequences and tripeptides in protein 3C

271    within a species was not detected (Figure 3). We can state that the tripeptide and motif KFRDI

272    are almost not variable within a species compared to domain d loop sequence, but there is a

273    specific preferred tripeptide sequence for each species. Hence, tripeptide sequences are species-

274    specific, while the domain d loop sequences are almost universal among species *Enterovirus A,*

275    *B, C, D* and *Rhinovirus A* and *C.*

<h1 style="text-align:center">Discussion</h1>

276

277       Most of the domain d apical loops in enterovirus genomes were represented by tetraloops.

278    The most common variants of tetraloop sequences corresponded to consensuses YNHG (Y=C/U,

279    N=any, H=A/C/U) (Table 2). Similar results were obtained in our previous experimental work,

280    where eight apical nucleotides of domain d of poliovirus genome were randomized, viable

281    variants were selected *in vitro* and the majority of selected tetraloops belonged to consensus

282    YNHG (Prostova *et al.*, 2015). Some tetraloops of consensus YNHG were found in genomes in

283    the NCBI database, but not among the variants selected *in vitro*, namely tetraloops CACG,

284    CUCG, UAAG, UGAG, CAAG, CGAG, UGUG, CUUG (Prostova *et al.*, 2015). Tetraloops

285    UGAG, UGUG and CUUG were reconstructed with U****G flanking base pair in context of

286    poliovirus genome strain Mahoney and supported effective virus replication (Prostova *et al.*,

287    2015).

288       Vice versa, tetraloops UUAG, UCAG, CCAG were found in domain d of selected *in vitro*

289    viable poliovirus variants and were able to support virus reproduction, but were not found in

290    naturally circulating viruses (Prostova *et al.*, 2015). One tetraloop of YNHG consensus (CUAG)

291    was not found neither in genomes from the NCBI database (n=2842), nor in the randomized

292    poliovirus genomes selected *in vitro* (n=62) (Table 2). Thermodynamic stability is unlikely to be

293    the reason why this and other tetraloops were unrepresented, as the melting temperature of stem

294    loops with avoided tetraloops is within range of the melting temperature of YNHG tetraloops

295    that supported replication (Proctor *et al.*, 2002). Moreover, tetraloops UUAG and UCAG are

296    common in rRNA (Woese *et al.*, 1990). Sample insufficiency cannot be excluded for both

297    database and *in vitro* selected sets of genomes, but it is safe to conclude that these tetraloop

298    variants are at least extremely rare. In any case, the fact that incidence of these tetraloops is

299    much less, than of other tetraloops points out that such variants are possibly less fit.

300            The most abundant tetraloops in the domain apical region of genomes from NCBI

301    database and variants selected *in vitro* could be compiled into consensuses UNCG and CNCG

302    (Table 2). At the same time, these tetraloops are the most abundant in rRNA, and with certain

303    closing base pairs are among the most thermodynamically stable tetraloops (Woese *et al.*, 1990;

304    Proctor *et al.*, 2002). Tetraloops of these consensuses and some other found tetraloops of YNHG

305    consensus form specific spatial structure of UNCG structural class of stable tetraloops (Cheong,

306    Varani and Tinoco, 1990; Varani, Cheong and Tinoco, 1991; Du *et al.*, 2003, 2004).

307            Another set of tetraloops, which correspond to GNYA consensus, was found both in

308    genomes of *Enterovirus E* and *F* and in genomes of viable polioviruses selected *in vitro*

309    (Prostova *et al.*, 2015). Tetraloop GCUA was able to support effective replication of poliovirus

310    and, together with tetraloop GUUA, are known to assume UNCG fold (Ihle *et al.*, 2005;

311    Melchers *et al.*, 2006; Prostova *et al.*, 2015).

312            Together, this data indicates that the spatial structure, rather than the exact sequence, is

313    the main subject for recognition by virus protein 3C, and, together with literature data, let us

314    assume that the sequence-structure degeneracy is an universal way RNA tetraloops are used in

315    nature (Lebars *et al.*, 2001; Wu *et al.*, 2004; Ihle *et al.*, 2005; D'Ascenzo *et al.*, 2016, 2017;

316    Bottaro and Lindorff-Larsen, 2017).

317         It can be speculated that pentaloops in domain d of *Enterovirus E* genome and triloops of

318    domain d of rhinoviruses have a potential to have the same UNCG fold as some YNHG and

319    GNYA tetraloops do. For HRV14 domain d it was shown that its triloop resembles the structure

320    of the first and two last nucleotides of UNCG structural class tetraloops (Headey *et al.*, 2007).

321    There are pentaloops with four nucleotides that belong to consensuses UNCG, GNRA or

322    gCUUGc and are able to form spatial structures of corresponding structural classes with the fifth

323    nucleotide bulged (Cai *et al.*, 1998; Schärpf *et al.*, 2000; Theimer, Finger and Feigon, 2003;

324    Oberstrass *et al.*, 2006; Liu *et al.*, 2009). It is possible that four nucleotides of pentaloops in

325    domain d of Enterovirus E species have UNCG fold with one bulged nucleotide.

326         Tetraloops that didn't belong to YNHG or GNYA consensuses were found in both sets of

327    natural and *in vitro* selected genomes. However, in an experiment such variants evolved towards

328    YNHG or GNYA consensuses (Prostova *et al.*, 2015). Apparently, tetraloops that don't belong to

329    YNHG or GNYA consensuses are less fit in most settings and in experimental conditions.

330    However, as these variants still could be found in few naturally circulating viruses

331    (consequently, they have emerged and have been fixed), we speculate that they may be beneficial

332    under specific replication conditions.

Peer Preprints

333      Similar structure of domain d and its apical region suggests free exchange of this region

334    between genomes of the same and of different species of Enterovirus genera. Indeed, viable intra

335    and inter species recombinants for this region could be obtained *in vitro* (Muslin *et al.*, 2015;

336    Bessaud *et al.*, 2016). To evaluate relative impact of the high mutation rate and recombination on

337    domain d apical loop variability, sequences of EV71 C4 genotype viruses were analyzed. Natural

338    recombination in EV71 genotype C4 is much less frequent than other EV-A types (Lukashev *et*

339    *al.*, 2014), and only one recombinant genome (accession number HQ423143) was detected in our

340    data set. Therefore variability of its domain d loop sequence reflects changes that were

341    accumulated via mutations only. Diversity of domain d loop sequence of EV-71 C4 viruses was

342    far less than among EV-A genomes and was represented only by 5 tetraloop sequence variants

343    (Table S1). The most recent common ancestor of EV71 genotype C4 dates about 20 years back

344    (McWilliam Leitch *et al.*, 2012), therefore this diversity, although limited, emerged very

345    recently. On the other hand, high sequence variability of domain d apical region in all

346    enterovirus genomes was possibly assisted by inter- and intraspecies recombination events.

347        Interestingly, in contrast to similar structure of domain d and very similar distribution of

348    its apical sequences in genomes of different enterovirus species, its putative RNA-recognition

349    tripeptide of 3C is diverse (Figure 3). Most of *Enterovirus A* genomes contain tripeptide VGK in

350    3C, while genomes of *Enterovirus B, C* and *D* species prevalently have TGK tripeptide (Figure

351    3). Genomes of *Rhinovirus A* and *C* also contain common enterovirus tetraloops in the domain d

352    apical region, but in the 3C they, unlike other species, contain tripeptides without positively

353    charged amino acids (Figure 3, Table S11, Table S12). Positively charged amino acids are often

354    involved into interaction with RNA, in particular with phosphates of the RNA backbone, thus

355    being of importance for RNA-protein recognition (Jones *et al.*, 2001; Bahadur, Zacharias and

356    Janin, 2008). In *Rhinovirus A* genomes positively charged amino acid "jumped" from the last

357    position of tripeptide (position 156) to the position that precedes tripeptide (position 153) (Figure

358    3, showed by arrow). Residue at position 153 starts and residue at position 156 ends reverse turn

359    between beta strands dII and eII of the protein 3C (Mosimann *et al.*, 1997; Matthews *et al.*, 1999;

360    Cui *et al.*, 2011).  In a crystal structure of Rhinovirus A2 protein 3C side chain of Lys153

361    (preceding tripeptide) is positioned in the similar region, as side chain of Lys156 (at the last

362    position of tripeptide) in crystal structure of Enterovirus 71 and Poliovirus 1 proteins 3C

363    (Mosimann *et al.*, 1997; Matthews *et al.*, 1999; Cui *et al.*, 2011). Thus, Lys at position 153 of 3C

364    has almost the same potential to interact with RNA-ligand as Lys at position 156  (Mosimann *et*

365    *al.*, 1997; Matthews *et al.*, 1999; Cui *et al.*, 2011). Genomes of *Rhinovirus C* species do not

366    contain a positively charged amino acid neither inside the tripeptide of 3C protein, nor in the

367    neighboring positions, possibly indicating that tripeptide 154-156 in protein 3C of *Rhinovirus C*

368    genome does not interact directly with RNA. Thereby, 3C is able to recognize domain d of oriL

369    with tripeptides of different sequence. In contrast to domain d structure and its apical sequence,

370    tripeptide is species-specific. Diversity of the tripeptide that is expected to recognize domain d

371    has several compatible explanations. The residue 154 of tripeptide possibly does not interact with

372    domain d directly. The tripeptide may be involved into a species-specific cooperative amino acid

373    network (amino acid "epistasis"). Moreover, different tripeptides might reflect slightly different

374    molecular mechanisms for domain d recognition.

375        The complexity of tripeptide's role in domain d recognition can be shown in several

376    examples. The 3C protein of different species with the same RNA-binding tripeptide is not

377    guaranteed to bind the same-structured domain d. Genomes of *Rhinovirus B* contain triloops in

378    the apical region of domain d that are paired with tripeptide TGK in 3C, common for genomes

379    with tetraloops. In contrast, protein 3C of Coxsackie virus B3 (*Enterovirus B* species, contain

380    tripeptide TGK) cannot recognize sufficiently well oriL with domain d capped with a triloop

381    (Zell *et al.*, 2002). This indicates that sequence of RNA-binding tripeptide probably is not the

382    exclusive participant in oriL-3C recognition. In other words, different molecular mechanisms of

383    oriL-3C recognition evolved in every enterovirus species independently. For example, it was

384    shown for Rhinovirus 14 (*Rhinovirus B* species) that protein 3C recognizes stem region of

385    domain d, rather than its apical loop (Leong *et al.*, 1993). Another oriL-3C recognition

386    mechanism apparently is used by *Enterovirus E* and *F* species, two oriLs of which play the same

387    role as the single oriL in genomes of other enteroviruses (Pilipenko, Blinov and Agol, 1990; Zell

388    *et al.*, 1999).  Apical loop of their domain d is tetra- or pentaloop with sequence that differs from

389    the loop consensuses of other enteroviruses. RNA-binding tripeptide in the 3C is species specific

390    as well, and is always MGK (Table S6). Interestingly, one genome of *Enterovirus G* species had

391    the same pair domain d loop – tripeptide of 3C, i.e. GUUA – MGK. Domain d of *Enterovirus G*

392    species is prolonged in comparison to length of domain d in genomes of other species

393    (Krumbholz *et al.*, 2002) (Figure 2). The tripeptide MGK in 3C of *Enterovirus E, F* and *G*

394    possibly indicates another molecular mechanism of oriL-3C recognition (Krumbholz *et al.*,

395    2002). Therefore, we assume that though putative RNA-binding tripeptide in most cases possibly

396    interacts with the domain d apical region (since amino acid substitutions in it are known to

397    compensate for structural disturbance in domain d), this interaction is not the only one that

398    determines the evolution oriL-3C interaction. Altogether, data demonstrates independent

399    evolution of putative RNA-binding tripeptide of 3C and domain d of oriL.

400

## Conclusions

402    We performed analysis of variety and occurrence of replication element oriL functional

403    loop and its protein ligand virus protease 3C. RNA-binding motifs of 3C are species-specific in

404    contrast to domain d loop sequences: domain d loop sequence variety is almost the same for

405    species *Enterovirus A, B, C, D* and *Rhinovirus A* and *C*, whereas tripeptide sequence variety

406    differ.

# References

407

408    Acevedo, A., Brodsky, L. and Andino, R. (2013) 'Mutational and fitness landscapes of an

409    RNA virus revealed through population sequencing.', *Nature*. Nature Publishing Group,

410    505(7485), pp. 686–690. doi: 10.1038/nature12861.

411    Andino, R., Rieckhof, G. E., Achacoso, P. L. and Baltimore, D. (1993) 'Poliovirus RNA

412    synthesis utilizes an RNP complex formed around the 5' -end of viral RNA', *EMBO Journal*,

413    12(9), pp. 3587–3598.

414    Andino, R., Rieckhof, G. E. and Baltimore, D. (1990) 'A Functional Ribonucleoprotein

415    around the 5' End of Poliovirus', *Cell*, 63, pp. 369–380.

416    Andino, R., Rieckhof, G. E., Trono, D. and Baltimore, D. (1990) 'Substitutions in the

417    protease (3Cpro) gene of poliovirus can suppress a mutation in the 5' noncoding region.', *J

418    Virol*, 64(2), pp. 607–612.

419    Bahadur, R. P., Zacharias, M. and Janin, J. (2008) 'Dissecting protein-RNA recognition

420    sites.', *Nucleic Acids Res.*, 36(8), pp. 2705–2716. doi: 10.1093/nar/gkn102.

421    Bessaud, M., Joffret, M.-L., Blondel, B. and Delpeyroux, F. (2016) 'Exchanges of

422    genomic domains between poliovirus and other cocirculating species C enteroviruses reveal a

423    high degree of plasticity', *Scientific Reports*, 6(1), p. 38831. doi: 10.1038/srep38831.

424    Bottaro, S. and Lindorff-Larsen, K. (2017) 'Mapping the Universe of RNA Tetraloop

425    Folds', *Biophysical Journal*, 113(2), pp. 257–267. doi: 10.1016/j.bpj.2017.06.011.

426      Cai, Z., Gorin, A., Frederick, R., Ye, X., Hu, W., Majumdar, A., Kettani, A. and Patel, D.

427    J. (1998) 'Solution structure of P22 transcriptional antitermination N peptide-box B RNA

428    complex.', *Nature structural biology*, 5(3), pp. 203–212.

429      Chase, A. J., Daijogo, S. and Semler, B. L. (2014) 'Inhibition of poliovirus-induced

430    cleavage of cellular protein PCBP2 reduces the levels of viral RNA replication.', *Journal of*

431    *virology*, 88(6), pp. 3192–3201. doi: 10.1128/JVI.02503-13.

432      Cheong, C. and Cheong, H. (2010) 'RNA Structure: Tetraloops', in *Encyclopedia of life*

433    *sciences*. Chichester: John Wiley & Sons, Ltd. doi: 10.1002/9780470015902.a0003135.pub2.

434      Cheong, C., Varani, G. and Tinoco, I. J. (1990) 'Solution structure of an unusually stable

435    RNA hairpin, 5'GGAC(UUCG)GUCC.', *Nature.*, 346(6285), pp. 680–682.

436      Crooks, G. E., Hon, G., Chandonia, J.-M. and Brenner, S. E. (2004) 'WebLogo: A

437    Sequence Logo Generator', *Genome Research*, 14(6), pp. 1188–1190. doi: 10.1101/gr.849004.

438      Cui, S., Wang, J., Fan, T., Qin, B., Guo, L., Lei, X., Wang, J., Wang, M. and Jin, Q.

439    (2011) 'Crystal structure of human enterovirus 71 3C protease.', *Journal of molecular biology*,

440    408(3), pp. 449–461. doi: 10.1016/j.jmb.2011.03.007.

441      D'Ascenzo, L., Leonarski, F., Vicens, Q. and Auffinger, P. (2016) '"Z-DNA like"

442    fragments in RNA: a recurring structural motif with implications for folding, RNA/protein

443    recognition and immune response.', *Nucleic acids research*, 44(12), pp. 5944–5956. doi:

444    10.1093/nar/gkw388.

445    D'Ascenzo, L., Leonarski, F., Vicens, Q. and Auffinger, P. (2017) 'Revisiting GNRA and

446    UNCG folds: U-turns versus Z-turns in RNA hairpin loops', *RNA*, 23(3), pp. 259–269. doi:

447    10.1261/rna.059097.116.

448    Du, Z., Yu, J., Andino, R. and James, T. L. (2003) 'Extending the family of UNCG-like

449    tetraloop motifs: NMR structure of a CACG tetraloop from coxsackievirus B3.', *Biochemistry*,

450    42(15), pp. 4373–4383. doi: 10.1021/bi027314e.

451    Du, Z., Yu, J., Ulyanov, N. B., Andino, R. and James, T. L. (2004) 'Solution structure of

452    a consensus stem-loop D RNA domain that plays important roles in regulating translation and

453    replication in enteroviruses and rhinoviruses.', *Biochemistry*, 43(38), pp. 11959–11972. doi:

454    10.1021/bi048973p.

455    Gamarnik, A. V and Andino, R. (1998) 'Switch from translation to RNA replication in a

456    positive-stranded RNA virus', *Genes & Dev.*, 12, pp. 2293–2304. doi: 10.1101/gad.12.15.2293.

457    Goodfellow, I., Chaudhry, Y., Richardson, A., Meredith, J., Almond, J. W., Barclay, W.

458    and Evans, D. J. (2000) 'Identification of a cis-acting replication element within the poliovirus

459    coding region.', *Journal of virology*, 74(10), pp. 4590–4600. doi: 10.1128/JVI.74.10.4590-

460    4600.2000.

461    Gruber, A. R., Lorenz, R., Bernhart, S. H., Neuböck, R. and Hofacker, I. L. (2008) 'The

462    Vienna RNA websuite.', *Nucleic acids research*, 36(Web Server issue), pp. W70-4. doi:

463    10.1093/nar/gkn188.

464    Hämmerle, T., Hellen, C. U. and Wimmer, E. (1991) 'Site-directed mutagenesis of the

465    putative catalytic triad of poliovirus 3C proteinase.', *The Journal of biological chemistry*, 266(9),

466    pp. 5412–5416.

467    Hämmerle, T., Molla, A. and Wimmer, E. (1992) 'Mutational analysis of the proposed

468    FG loop of poliovirus proteinase 3C identifies amino acids that are necessary for 3CD cleavage

469    and might be determinants of a function distinct from proteolytic activity.', *J Virol.*, 66(10), pp.

470    6028–6034.

471    Harris, K. S., Reddigari, S. R., Nicklin, M. J., Hämmerle, T. and Wimmer, E. (1992)

472    'Purification and characterization of poliovirus polypeptide 3CD, a proteinase and a precursor

473    for RNA polymerase.', *J Virol.*, 66(12), pp. 7481–7489.

474    Headey, S. J., Huang, H., Claridge, J. K., Soares, G. A., Dutta, K., Schwalbe, M., Yang,

475    D. and Pascal, S. M. (2007) 'NMR structure of stem-loop D from human rhinovirus-14.', *RNA*,

476    13(3), pp. 351–360. doi: 10.1261/rna.313707.

477    Ihle, Y., Ohlenschläger, O., Häfner, S., Duchardt, E., Zacharias, M., Seitz, S., Zell, R.,

478    Ramachandran, R. and Görlach, M. (2005) 'A novel cGUUAg tetraloop structure with a

479    conserved yYNMGg-type backbone conformation from cloverleaf 1 of bovine enterovirus 1

480    RNA.', *Nucleic Acids Res.*, 33(6), pp. 2003–2011. doi: 10.1093/nar/gki501.

481    Jones, S., Daley, D. T. A., Luscombe, N. M., Berman, H. M. and Thornton, J. M. (2001)

482    'Protein – RNA interactions : a structural analysis', *Biochemistry*, 29(4), pp. 943–954.

483     Krumbholz, A., Dauber, M., Henke, A., Birch-Hirschfeld, E., Knowles, N. J., Stelzner, A.

484     and Zell, R. (2002) 'Sequencing of porcine enterovirus groups II and III reveals unique features

485     of both virus groups.', *Journal of virology*, 76(11), pp. 5813–5821.

486     Larkin, M. A., Blackshields, G., Brown, N. P., Chenna, R., McGettigan, P. A.,

487     McWilliam, H., Valentin, F., Wallace, I. M., Wilm, A., Lopez, R., Thompson, J. D., Gibson, T.

488     J. and Higgins, D. G. (2007) 'Clustal W and Clustal X version 2.0.', *Bioinformatics*, 23(21), pp.

489     2947–2948. doi: 10.1093/bioinformatics/btm404.

490     Lauring, A. S., Frydman, J. and Andino, R. (2013) 'The role of mutational robustness in

491     RNA virus evolution.', *Nature reviews. Microbiology*, 11(5), pp. 327–336. doi:

492     10.1038/nrmicro3003.

493     Lebars, I., Lamontagne, B., Yoshizawa, S., Aboul-Elela, S. and Fourmy, D. (2001)

494     'Solution structure of conserved AGNN tetraloops: insights into Rnt1p RNA processing.', *The*

495     *EMBO journal*, 20(24), pp. 7250–7258. doi: 10.1093/emboj/20.24.7250.

496     Leong, L. E. C., Walker, P. A., Porter, A. G., Protease, H. R.-, Leon, L. E. C., Walker, P.

497     A., Porter, A. G., Leong, L. E. C., Walker, P. A. and Porter, A. G. (1993) 'Human Rhinovirus-14

498     Protease 3C (3Cpro) Binds Specifically to the 5'-Noncoding Region of the Viral RNA', *The*

499     *Journal of biological chemistry*, 268(34), pp. 25735–25739.

500     Liu, P., Li, L., Keane, S. C., Yang, D., Leibowitz, J. L. and Giedroc, D. P. (2009) 'Mouse

501     hepatitis virus stem-loop 2 adopts a uYNMG(U)a-like tetraloop structure that is highly

502    functionally tolerant of base substitutions.', *Journal of virology*, 83(23), pp. 12084–12093. doi:

503    10.1128/JVI.00915-09.

504          Lorenz, R., Bernhart, S. H., Höner Zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P.

505    F. and Hofacker, I. L. (2011) 'ViennaRNA Package 2.0.', *Algorithms for molecular biology*,

506    6(26). doi: 10.1186/1748-7188-6-26.

507          Lukashev, A. N., Shumilina, E. Y., Belalov, I. S., Ivanova, O. E., Eremeeva, T. P.,

508    Reznik, V. I., Trotsenko, O. E., Drexler, J. F. and Drosten, C. (2014) 'Recombination strategies

509    and evolutionary dynamics of the Human enterovirus A global gene pool.', *The Journal of*

510    *general virology*, 95(Pt 4), pp. 868–873. doi: 10.1099/vir.0.060004-0.

511          Matthews, D. a, Dragovich, P. S., Webber, S. E., Fuhrman, S. a, Patick,  a K., Zalman, L.

512    S., Hendrickson, T. F., Love, R. a, Prins, T. J., Marakovits, J. T., Zhou, R., Tikhe, J., Ford, C. E.,

513    Meador, J. W., Ferre, R. a, Brown, E. L., Binford, S. L., Brothers, M. a, DeLisle, D. M. and

514    Worland, S. T. (1999) 'Structure-assisted design of mechanism-based irreversible inhibitors of

515    human rhinovirus 3C protease with potent antiviral activity against multiple rhinovirus

516    serotypes.', *Proceedings of the National Academy of Sciences of the United States of America*,

517    96(20), pp. 11000–11007.

518          McWilliam Leitch, E. C., Cabrerizo, M., Cardosa, J., Harvala, H., Ivanova, O. E., Koike,

519    S., Kroes, A. C. M., Lukashev, A., Perera, D., Roivainen, M., Susi, P., Trallero, G., Evans, D. J.

520    and Simmonds, P. (2012) 'The association of recombination events in the founding and

521    emergence of subgenogroup evolutionary lineages of human enterovirus 71.', *Journal of*

522    *virology*, 86(5), pp. 2676–2685. doi: 10.1128/JVI.06065-11.

523          Melchers, W. J. G., Zoll, J., Tessari, M., Bakhmutov, D. V, Gmyl, A. P., Agol, V. I. and

524    Heus, H. a (2006) 'A GCUA tetranucleotide loop found in the poliovirus oriL by in vivo SELEX

525    (un)expectedly forms a YNMG-like structure: Extending the YNMG family with GYYA.', *RNA*,

526    12(9), pp. 1671–1682. doi: 10.1261/rna.113106.

527          Mosimann, S. C., Cherney, M. M., Sia, S., Plotch, S. and James, M. N. (1997) 'Refined

528    X-ray crystallographic structure of the poliovirus 3C gene product.', *J Mol Biol*, 273(5), pp.

529    1032–1047. doi: 10.1006/jmbi.1997.1306.

530          Muslin, C., Joffret, M.-L., Pelletier, I., Blondel, B. and Delpeyroux, F. (2015) 'Evolution

531    and Emergence of Enteroviruses through Intra- and Inter-species Recombination: Plasticity and

532    Phenotypic Impact of Modular Genetic Exchanges in the 5' Untranslated Region', *PLOS*

533    *Pathogens*, 11(11), p. e1005266. doi: 10.1371/journal.ppat.1005266.

534          Oberstrass, F. C., Lee, A., Stefl, R., Janis, M., Chanfreau, G. and Allain, F. H.-T. (2006)

535    'Shape-specific recognition in the structure of the Vts1p SAM domain with RNA.', *Nature*

536    *structural & molecular biology*, 13(2), pp. 160–167. doi: 10.1038/nsmb1038.

537          Oermann, C. M., Schuster, J. E., Conners, G. P., Newland, J. G., Selvarangan, R. and

538    Jackson, M. A. (2015) 'Enterovirus D68. A focused review and clinical highlights from the 2014

539    U.S. Outbreak.', *Annals of the American Thoracic Society*, 12(5), pp. 775–781. doi:

540    10.1513/AnnalsATS.201412-592FR.

541        Palmenberg, A., Neubauer, D. and Skern, T. (2010) 'Genome organization and encoded

542    proteins.', in Ehrenfeld, E., Domingo, E., and Roos, R. P. (eds) *The Picornaviruses*. ASM Press,

543    pp. 3–17.

544        Pilipenko, E. V, Blinov, V. M. and Agol, V. I. (1990) 'Gross rearrangements within the

545    5'-untranslated region of the picornaviral genomes.', *Nucleic acids research*, 18(11), pp. 3371–

546    3375.

547        Proctor, D. J., Schaak, J. E., Bevilacqua, J. M., Falzone, C. J. and Bevilacqua, P. C.

548    (2002) 'Isolation and characterization of a family of stable RNA tetraloops with the motif

549    YNMG that participate in tertiary interactions.', *Biochemistry*, 41(40), pp. 12062–12075.

550        Prostova, M. A., Gmyl, A. P., Bakhmutov, D. V, Shishova, A. A., Khitrina, E. V,

551    Kolesnikova, M. S., Serebryakova, M. V, Isaeva, O. V and Agol, V. I. (2015) 'Mutational

552    robustness and resilience of a replicative cis-element of RNA virus: promiscuity, limitations,

553    relevance.', *RNA biology*, 12(12), pp. 1338–1354. doi: 10.1080/15476286.2015.1100794.

554        Rieder, E., Xiang, W., Paul, A. and Wimmer, E. (2003) 'Analysis of the cloverleaf

555    element in a human rhinovirus type 14/poliovirus chimera: correlation of subdomain D structure,

556    ternary protein complex formation and virus replication', *J Gen Virol*, 84(8), pp. 2203–2216. doi:

557    10.1099/vir.0.19013-0.

558        Sanjuán, R., Nebot, M. R., Chirico, N., Mansky, L. M. and Belshaw, R. (2010) 'Viral

559     mutation rates.', *Journal of virology*, 84(19), pp. 9733–9748. doi: 10.1128/JVI.00694-10.

560     Schärpf, M., Sticht, H., Schweimer, K., Boehm, M., Hoffmann, S. and Rösch, P. (2000)

561     'Antitermination in bacteriophage lambda. The structure of the N36 peptide-boxB RNA

562     complex.', *Eur. J. Boichem.*, 267(267), pp. 2397–2408.

563     Shih, S., Chen, T. and Wu, C. (2004) 'Mutations at KFRDI and VGK Domains of

564     Enterovirus 71 3C Protease Affect Its RNA Binding and Proteolytic Activities', *Journal of*

565     *Biomedical Science*, pp. 239–248. doi: 10.1159/000076036.

566     Solomon, T., Lewthwaite, P., Perera, D., Cardosa, M. J., McMinn, P. and Ooi, M. H.

567     (2010) 'Virology, epidemiology, pathogenesis, and control of enterovirus 71.', *The Lancet.*

568     *Infectious diseases*, 10(11), pp. 778–790. doi: 10.1016/S1473-3099(10)70194-8.

569     Theimer, C. A., Finger, L. D. and Feigon, J. (2003) 'YNMG tetraloop formation by a

570     dyskeratosis congenita mutation in human telomerase RNA', *RNA*, 9, pp. 1446–1455. doi:

571     10.1261/rna.5152303.activity.

572     Thompson, A. A. and Peersen, O. B. (2004) 'Structural basis for proteolysis-dependent

573     activation of the poliovirus RNA-dependent RNA polymerase.', *The EMBO journal*, 23(17), pp.

574     3462–3471. doi: 10.1038/sj.emboj.7600357.

575     Trono, D., Andino, R. and Baltimore, D. (1988) 'An RNA sequence of hundreds of

576     nucleotides at the 5' end of poliovirus RNA is involved in allowing viral protein synthesis.',

577     *Journal of virology*, 62(7), pp. 2291–2299.

578     Uhlenbeck, O. C. (1990) 'Tetraloops and RNA folding', *Nature*. Nature Publishing

579   Group, 346(6285), pp. 613–614. doi: 10.1038/346613a0.

580     Varani, G., Cheong, C. and Tinoco, I. (1991) 'Structure of an unusually stable RNA

581   hairpin.', *Biochemistry*, 30(13), pp. 3280–3289.

582     Vogt, D. A. and Andino, R. (2010) 'An RNA element at the 5'-end of the poliovirus

583   genome functions as a general promoter for RNA synthesis', *PLoS pathogens*, 6(6), p. e1000936.

584   doi: 10.1371/journal.ppat.1000936.

585     Wagner, A. and Stadler, P. F. (1999) 'Viral RNA and evolved mutational robustness.',

586   *The Journal of experimental zoology*, 285(2), pp. 119–127.

587     Woese, C. R., Winker, S., Gutell, R. R., Winkers, S. and Gutell, R. R. (1990)

588   'Architecture of ribosomal RNA : Constraints on the sequence of "tetra-loops"', *Proceedings of*

589   *the National Academy of Sciences of the United States of America*, 87(November), pp. 8467–

590   8471.

591     Wu, H., Henras, A., Chanfreau, G. and Feigon, J. (2004) 'Structural basis for recognition

592   of the AGNN tetraloop RNA fold by the double-stranded RNA-binding domain of Rnt1p RNase

593   III', *Amino Acids*, 101(22), pp. 8307–8312.

594     Zell, R., Sidigi, K., Bucci, E., Stelzner, A. and Görlach, M. (2002) 'Determinants of the

595   recognition of enteroviral cloverleaf RNA by coxsackievirus B3 proteinase 3C.', *RNA (New*

596   *York, N.Y.)*, 8(2), pp. 188–201. Available at:

597     http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1370242&tool=pmcentrez&renderty

598     pe=abstract (Accessed: 19 December 2014).

599          Zell, R., Sidigi, K., Henke, A., Schmidt-Brauns, J., Hoey, E., Martin, S. and Stelzner, A.

600     (1999) 'Functional features of the bovine enterovirus 5'-non-translated region.', *J Gen Virol*, 80,

601     pp. 2299–2309.

602

603

# Figure 1

Schematic representation of poliovirus genome and detailed representation of secondary structure of poliovirus replicative element oriL.
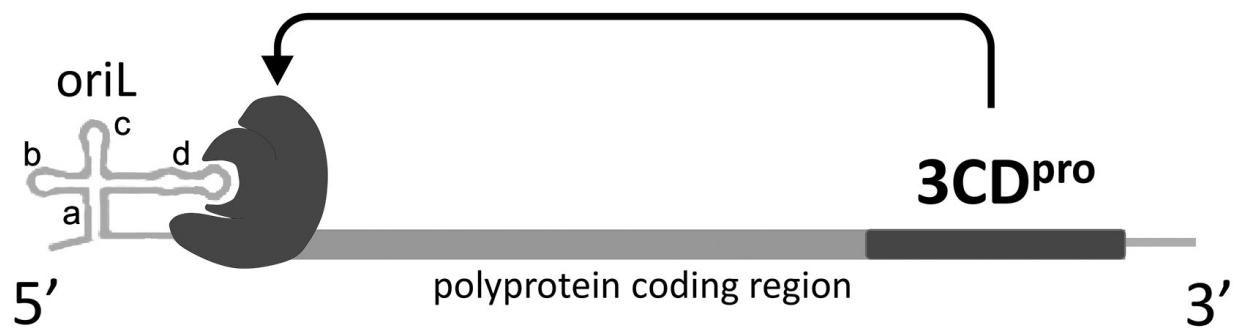
# Figure 2

Secondary structure of oriL domain d of distinct enterovirus species.

For *Enterovirus E* and *F* domain d of the first oriL is shown. Secondary structure of domain d of Porcine enterovirus 9 strain UKG/410/73 was folded with use as reference Krumbholtz et al., 2002.
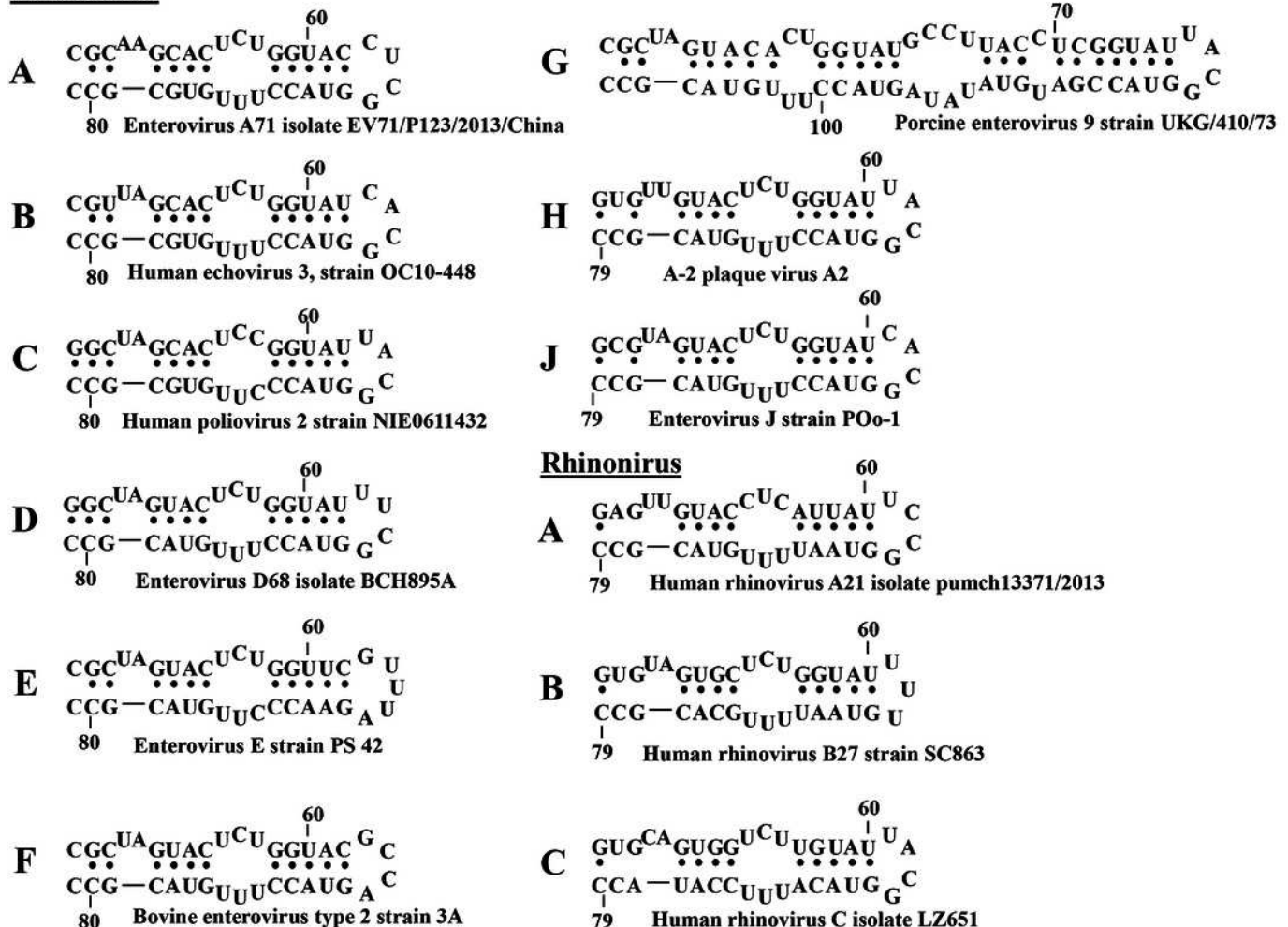
# Figure 3

Distribituion of domain d loop sequence and amino acid motifs in the 3C protein.

**A -** Distribution of domain d loop sequences. The regions corresponding to tetraloop consensuses, triloops and pentaloops are shown. Number of genomes cut off at 15 for clear view of sequence distribution. **B -** The frequency plot of amino acid sequence of 3C in species of genus *Enterovirus*. The amino acid sequence logo was done with WebLogo server (Crooks *et al.*, 2004). Arrows indicate amino acids of the proteolytic triade (Glu71 and Cys 147), the first and the last amino acids of motif $_{82}$KFRDI$_{86}$, the putative RNA-binding tripeptide 154-156 of 3C and Lys153 in the protein 3C of *Rhinovirus A*.
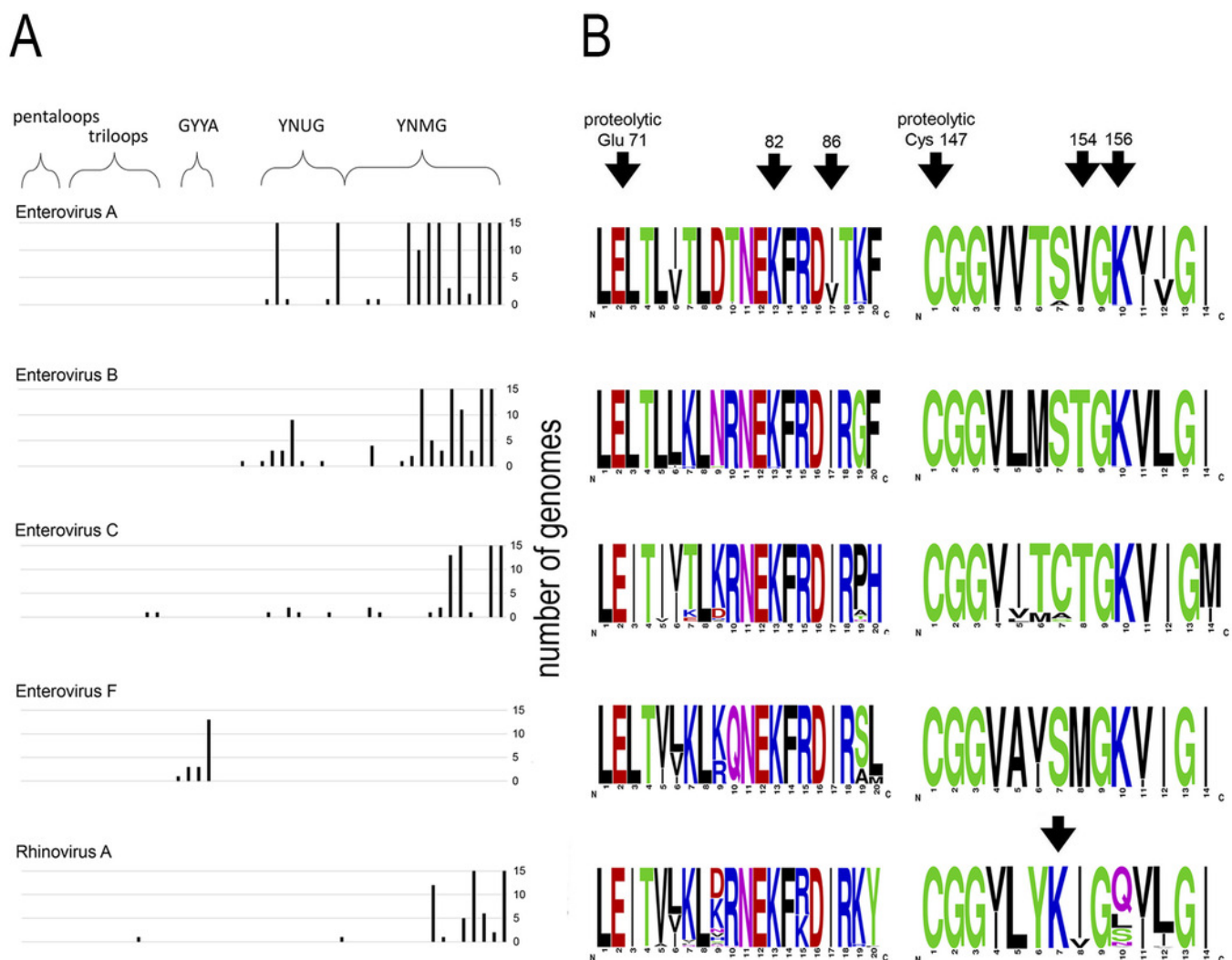
**Table 1**(on next page)

Number of full genome sequences that contained oriL region and number of unique domain d sequences before and after filtration.

For *Enterovirus E* and *F* number of unique tetraloops is shown separately for first and the second oriL.

| Species | Number of full genome sequences | Number of full genome sequences after 1% nucleic identity filtration | Number of unique tetraloops before filtration | | Number of unique tetraloops after filtration | |
|---|---|---|---|---|---|---|
| Enterovirus A | 1052 | 564 | 17 | | 16 | |
| Enterovirus B | 339 | 244 | 18 | | 18 | |
| Enterovirus C | 747 | 274 | 15 | | 12 | |
| Enterovirus D | 419 | 57 | 7 | | 6 | |
| Enterovirus E | 12 | 10 | 6 | 5 | 6 | 5 |
| Enterovirus F | 13 | 10 | 4 | 3 | 4 | 3 |
| Enterovirus G | 10 | 8 | 6 | | 6 | |
| Enterovirus H | 3 | 2 | 2 | | 2 | |
| Enterovirus J | 8 | 5 | 3 | | 3 | |
| Rhinovirus A | 151 | 118 | 8 | | 8 | |
| Rhinovirus B | 50 | 37 | 7 | | 7 | |
| Rhinovirus C | 38 | 37 | 6 | | 6 | |

1

**Table 2**(on next page)

Occurrence of domain d apical sequences in filtered sets of full genomes of different enterovirus species.

Tetraloops CCCG, UGUG, CAUG and UUGG that were unique for species *Enterovirus A, C* and *D* and were lost upon filtration, were added and are shown in blue. The gradient coloring from red to green represents abundancy heat map for the genomes with different domain d sequence.

| Loop sequence | Enterovirus | | | | | | | | | Rhinovirus | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | B | C | D | E | F | G | H | J | A | B | C |
| **Triloops** | | | | | | | | | | | | |
| CCG | | | 1 | | | | | | | | | |
| CAG | | | 1 | | | | | | | | | |
| UCU | | | | | | | | | | 1 | 5 | |
| UUU | | | | | | | | | | | 17 | |
| UAU | | | | | | | | | | | 8 | |
| AUU | | | | | | | | | | | 4 | |
| UGU | | | | | | | | | | | 1 | |
| UUC | | | | | | | | | | | 1 | |
| GAU | | | | | | | | | | | 1 | |
| **YNMG[a] Tetraloops** | | | | | | | | | | | | |
| UACG | 85 | 51 | 106 | | | | 3 | 1 | 2 | 38 | | 15 |
| UGCG | 114 | 31 | 43 | | | | 2 | 1 | | 2 | | |
| UUCG | 16 | 3 | | 50 | | | | | | 6 | | 6 |
| UCCG | 2 | 11 | 1 | | | | | | | 53 | | 10 |
| CACG | 48 | 98 | 101 | 1 | | | 1 | | 2 | 5 | | |
| CGCG | 3 | 3 | 13 | | | | 1 | | | | | |
| CUCG | 132 | 5 | 2 | 2 | | | | | | 1 | | 3 |
| CCCG | 40 | 16 | 1 | 1 | | | | | | 12 | | 2 |
| UAAG | 10 | 2 | | | | | | | | | | |
| UGAG | 22 | 1 | | | | | | | | | | |
| UUAG | | | | | | | | | | | | |
| UCAG | | | | | | | | | | | | |
| CAAG | 1 | 4 | 1 | | | | | | 1 | | | |
| CGAG | 1 | | 2 | | | | | | | | | |
| CUAG | | | | | | | | | | | | |
| CCAG | | | | | | | | | | | | |
| YACG | | 1 | | | | | | | | | | |
| **YNUG Tetraloops** | | | | | | | | | | | | |
| UAUG | 54 | | | | | | | | | | | |
| UGUG | 1 | 1 | 1 | | | | | | | | | |
| UUUG | | | | 1 | | | | | | | | |
| UCUG | | | 1 | | | | | | | | | |
| CAUG | | | 9 | 1 | | | | | | | | |
| CGUG | 1 | 3 | 2 | | | | | | | | | |
| CUUG | 34 | 3 | | 2 | | | | | | | | |
| CCUG | 1 | 1 | 1 | | | | | | | | | |
| **GYYA Tetraloops** | | | | | | | | | | | | |
| GCUA | | | | | 2 | 13 | | | | | | |
| GCCA | | | | | | 3 | | | | | | |
| GUUA | | | | | 2 | 3 | 1 | | | | | |
| **Other tetraloops** | | | | | | | | | | | | |
| UUGG | | | | | 1 | | | | | | | |
| CUUC | | | | | | | | | | | | 1 |
| AUUA | | | | | | 1 | | | | | | |
| **Pentaloops** | | | | | | | | | | | | |
| GCUUA | | | | | 7 | | | | | | | |
| GUUUA | | | | | 2 | | | | | | | |
| GCCUA | | | | | 4 | | | | | | | |
| GCGUA | | | | | 1 | | | | | | | |
| GAUUA | | | | | 1 | | | | | | | |
| GUCUA | | | | | 1 | | | | | | | |

Legend: 1   2   11   31   51   132

1   a – here and after Y=C/U, N=any nucleotide, M=A,C