

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Analysis of the contribution of intrinsic disorder in shaping potyvirus genetic diversity.

Guillaume Lafforgue, Thierry Michon and Justine Charon

Abstract

Intrinsically disordered regions (IDRs) are abundant in the proteome of RNA viruses. The multifunctional properties of these regions are widely documented and their structural flexibility is associated with low constraint in the amino acid positions. Therefore, from an evolutionary stand point, these regions could have a greater mutational permissiveness than highly structured regions (ORs for Ordered Regions). They could thus provide a potential adaptive reservoir. To address this hypothesis, we compared the mutational robustness of IDRs and ORs in the genome of potyviruses, a major genus of plant viruses. For this purpose, a simulation model (DOI: 10.5281/zenodo.6396239) was built and used to distinguish a possible selection phenomenon in the biological data sets from randomly generated mutations. We analyzed several short-term experimental evolution datasets. An analysis was also performed on the natural diversity of three different species of potyviruses reflecting the long-term evolution. We observed that the mutational robustness of IDRs is significantly higher than that of ORs. Moreover, the substitutions in the ORs are very constrained by the conservation of the physico-chemical properties of the amino acids. This feature is not found in the IDRs where the substitutions tend to be more random. This reflects the weak structural constraints in these regions, in which an amino acid polymorphism is naturally conserved in the course of evolution, potyvirus IDRs and ODRs follow different evolutive paths with respect to their mutational robustness. These results force to consider the hypothesis that during selection, adaptive solutions could emerge from the amino acid polymorphism carried by IDRs.

Introduction

Protein intrinsic disorder

Proteins possess intrinsically disordered regions (IDRs), i.e. regions lacking a unique three dimensional structure and yet capable of exerting important biological functions [1,2], which challenges the so-called “structure-function relationship” dogma. Although it is today quite admitted that intrinsically disordered hub proteins are key players in the cellular interactome, the involvement of intrinsically disorder (ID) in evolution is still under debate. An earlier study was aimed at comparing the structural features of single-domain small- proteins from hypothermophilic bacteria, archaea, mesophilic eukaryota and prokaryota, and RNA or DNA viruses, whose crystal structures were available [3]. It was concluded from this analysis that viral proteins and more particularly RNA virus proteins, display (i) higher stability upon simulations of mutation accumulation and (ii) lower inter-residues contact densities. This latter feature is a typical signature of intrinsic disorder. It has thus been proposed that the large intrinsic disorder content in viral proteins could contribute to efficiently buffer mutation effects [3,4]. This was experimentally shown in the case of the intrinsically disordered protein VPg from potyviruses [5]. This is strongly contrasting with non-additive/epistatic stability loss profile expected from ordered proteins as previously reported for a bacterial β -lactamase [6]. It is hence conceivable that low structural requirements in IDRs could lead to some mutational robustness, and in turn, to an easier way for exploring the mutational space, without dramatic impairment of the protein biological functions. For instance, this idea sounds especially relevant regarding RNA virus adaptation to the host. This could contribute to a rapid adaptation to environmental stresses, without excessive loss of fitness. There is no doubt that this question is of very general interest. Consequently, the high evolutionary potential of RNA viruses, and the high ID content in their proteins, set the basis for assessing the

49 contribution of ID to the shaping of virus genetic diversity in a context of host adaptation. Plant-
50 phytovirus pathosystems provide useful experimental models for studying these aspects [7].
51 An *in silico* analysis unveiled a high ID content in the *Potyvirus* proteome both at inter- and intra-
52 species scales [8]. This feature has been conserved during *Potyvirus* evolution, suggesting a functional
53 advantage of ID. When comparing the evolutionary constraint (ratio of non-synonymous to
54 synonymous substitution rates, dN/dS) between ordered and disordered regions within the proteome
55 of different potyvirus species, IDRs display significantly higher dN/dS values than ordered regions
56 (ORs), a finding that indicates a tendency of intrinsically disordered domains to evolve faster than more
57 structured regions during potyvirus evolution [8]. Using the pathosystem PVY/pepper, we previously
58 obtained the first *in vivo* experimental data supporting the hypothesis that IDRs could influence virus
59 adaptability to the host [9], possibly by enabling a faster exploration of the mutational space, thereby
60 allowing the virus to bypass the plant resistance. Indeed, a correlation was observed between the
61 adaptive potential of the virus and the disorder content within the VPg viral protein.

62
63 To further assess this previously described role of IDRs on RNA virus adaptation, the present study
64 aimed at analyzing whether the regions predicted as disordered in viral proteomes are more likely to
65 evolve and accommodate amino acid substitutions (non-synonymous mutations) than more structured
66 areas. Ordered and disordered region sequences from various potyvirus species were thus retrieved
67 and compared for several adaptive parameters at two timescales of viral evolution, a short-term scale
68 experimental evolution and a long-term evolution reflected by natural diversity.

69 The short-term scale data analyzed consisted in high-throughput sequencing (HTS) retrieved from
70 three independent evolution experiments, i.e. PVY [10,11], and TEV [12]. HTS provides access to the
71 complete genome sequences of all viral variants - including those that are in a minority - that make up
72 a population [13]. By sequencing each individual genome from the viral population, it is thus possible
73 to assess the genetic structures of evolving potyvirus populations and thus potentially address the
74 processes that shape this genetic variability, and to a greater extent the evolvability of the viral
75 population.

76 To evaluate the impact of disordered versus ordered region on potyvirus evolvability (i.e. mutational
77 robustness) at a higher scale of evolution, genomic sequences from TuMV, TEV and PVY natural
78 diversity were also retrieved.

79 To prevent bias in our analysis of the structural determinant on potyvirus evolution, a third dataset,
80 corresponding to simulated data was also obtained. Briefly, potyvirus genomes were artificially
81 mutated *in silico* according the viral replicase features, to mimic the genetic diversity obtained in the
82 absence of selection and, among others, effects of protein structural determinants. Adaptive
83 parameters of the resulting mutants were thus obtained and compared to those from the biological
84 data.

85

86 **Material and methods**

87 **Data sets**

88 ○ Disorder prediction

89 We scanned for disordered regions along potyvirus polyproteins using Predictor of Naturally
90 Disordered Regions (PONDR-VLXT), an algorithm accessible through the Disprot server
91 (<http://disorder.compbio.iupui.edu/metapredictor.php>) [14,15]. Parameters were set to “default” for
92 ID score predictions.

93 ○ Experimental dataset

94 The study of Cuevas et al 2015 (TEV 2015), [12] evolved the TEV on two different host, *Nicotiana*
95 *tabacum* and *Capsicum annuum*, while the two others studies Kutnjak et al 2015, 2017 [10,11] (PVY
96 2015 and PVY 2017), used the PVY on *Solanum tuberosum*. Table S1 compiles all the resulting
97 mutations of experimental datasets.

98 ○ Natural diversity dataset

99 Datasets used contained 6 genomes of TEV isolates, 100 genomes of PVY isolates and 100 genomes
100 of TuMV isolates. Corresponding genome accessions are listed in Table S2. These datasets will be
101 referred to as TEV_{ND}, PVY_{ND} and TuMV_{ND} in the study.

102 ○ Simulation

103 The distribution of mutations in the virus sequence is the sum of the contribution of viral
104 polymerase errors and of the subsequent selection according to structure-function relationships. In
105 order to uncouple these two components, we built an algorithm to mimic the distribution of
106 synonymous and non-synonymous mutations introduced by the low fidelity virus RNA polymerase
107 during genome replication (DOI: 10.5281/zenodo.6396239). It was hypothesized that, mutations could
108 be randomly introduced all along the genome during its replication. Consequently, if IDRs and ORs
109 were equally susceptible to mutations, NS and S were expected to be homogeneously distributed in
110 each of the two regions before virus submission to the selection pressure. The simulation takes also
111 into consideration the specificity of viral polymerase on transversion/transition mutations calculated
112 from TEV experimental data [16].

113 We generated n variants from the original potyvirus sequence, with each variant bearing one SNP .

114

115 ○ Adaptive components tested

116 The collections of sequences generated from experimental evolution, natural diversity and
117 simulated experiments were then analyzed with respect to the number of S and NS mutations (DOI:
118 10.5281/zenodo.6396239) and for each viral protein, their location either in the ORs or IDRs. BLOSUM-
119 based scores of each NS mutations were also used to determine the potential of IDRs and ORs to cope
120 with amino acid substitutions (DOI: 10.5281/zenodo.6396239). Finally, the characteristic of naturally
121 occurring substitutions were analyzed in term of maintenance versus disturbance of disorder. It was
122 reported that ORs and IDRs possess distinct sequence biases. Promotor scores, ranging from 0 to 1 (1
123 being the highest promoter score for disorder) were adapted from a previously published classification
124 [17] and associated to each amino acids (Table S3).

125

126 **Results**

127 To evaluate the contribution of intrinsically disordered regions on potyvirus evolvability and
128 adaptation, this study compared viral genomic populations retrieved from evolution experiments,
129 natural diversity and *in silico* generated pool of variants. Several parameters were thus assessed and
130 compared at both genomic and proteomic levels and consisted in (i) the location of the diversity (within
131 intrinsically disordered versus ordered protein or regions), (ii) the nature of nucleotide mutations
132 (synonymous versus non-synonymous) as well as (iii) the biochemical and disorder-promoting nature
133 of corresponding amino acid substitutions.

134

135 **Theoretical minimum number of mutations required for an accurate estimation of S and NS**
136 **distribution in the genome**

137 Datasets generated from the three experimental evolutions represent between 115 and 317
138 mutations. We hypothesized that the number of mutations considered could be too low to lead to
139 robust conclusions. Consequently, a first consideration was to estimate the average number of
140 mutations required to be significant. Four independent generations of 100, 300, 500, 750, 1000 and
141 1250 mutations were thus randomly introduced along the TEV genome sequence. Assuming that such
142 random mutagenesis should not be impacted by any structural or protein determinants, the number
143 of mutations (synonymous and non-synonymous) should be equally distributed along the genome,
144 regardless of the corresponding proteome intrinsic disorder. Thus, the distribution of NS and S
145 mutations among IDRs were determined (Figure 1). Above 600 mutations, an equal distribution of
146 mutations among either ORs and IDRs was observed. Therefore, in order to ensure representative

147 values for further analysis, the results of 4 independent simulations with 1000 mutations each, will be
148 used.
149 This threshold of 1000 mutations, which is required for a robust analysis, was confirmed by monitoring
150 the evolution of the R^2 coefficient as a function of the mutation number. Whether for NS or S, below
151 750 mutations, the R^2 greatly fluctuates (Figure S1).

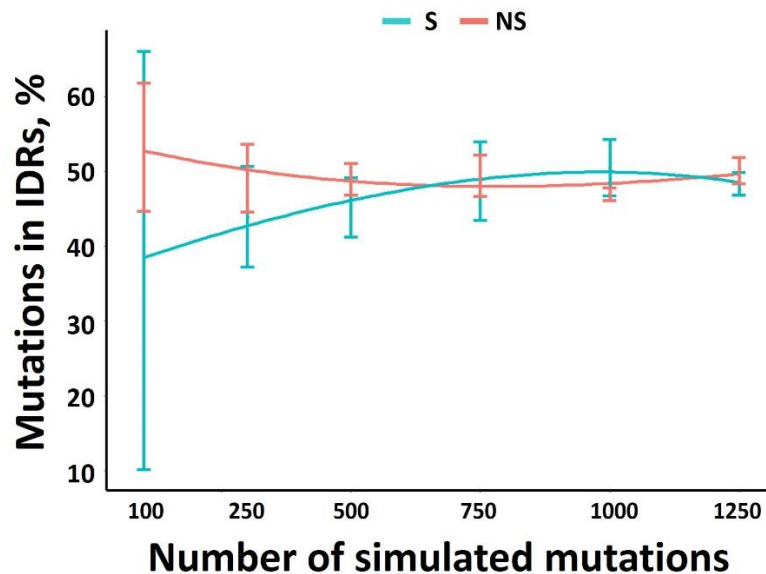


Figure 1. The % of S and NS mutations in IDRs versus the mutations number in the TEV genome. For a given number of mutations, 4 independent simulations were run.

152
153 This result confirms that the limited number of mutations available from the experimental data sets is
154 likely to make our analysis less robust. To increase the size of the dataset and extend our observations
155 to larger scales of viral evolution, the natural diversity of TEV_{ND}, PVY_{ND} and TuMV_{ND} isolates was also
156 analyzed by retrieving complete genomes available in Genbank. With 1296, 4646 and 7528 mutations
157 identified in the corresponding TEV_{ND}, PVY_{ND} and TuMV_{ND} datasets. These data should allow us to
158 assess whether there is a significant difference in mutational robustness between IDRs and ORs.

159 160 **Correlation assessment between protein length and number of mutations**

161 We first assessed the propensity of each potyvirus proteins to accumulate adaptive non-synonymous
162 (NS) versus synonymous (S) mutations. The number of NS or S mutations observed in each protein
163 coding sequence divided by the total protein length were thus calculated for each of the experimental
164 evolution, natural diversity and simulated data sets.

165 At the short-term evolution scale, the longer the protein, the higher the number of S mutations, with
166 a significant correlation between protein length and percentage of S mutations. By contrast NS
167 mutation number were not correlated with the protein length, for the three experimental studies
168 analyzed [10–12] (Table 1).

169 At the long-term evolution scale, the natural diversity confirmed the trend that the accumulation of
170 NS poorly correlated with protein length. Non-synonymous adaptive mutations, which reflect the viral
171 amino acid polymorphism, are thus not accumulated homogenously along the potyvirus proteome.

172
173
174
175
176
177
178

	S	NS
TEV 2015	0,63	0,19
TEV _{ND}	0,94	0,16
Simulations	1	0,98

	S	NS
PVY 2015	0,93	0,12
PVY 2017	0,78	0,01
PVY _{ND}	0,96	0,35
Simulations	0,96	0,97

	S	NS
TuMV _{ND}	0,95	0,09
Simulations	0,92	0,98

Table 1. Correlation coefficient (R^2) between coding sequence length of the TEV proteins and the mutations (S or NS). Experimental evolution: TEV 2015 [12], PVY 2015 [10] and PVY 2017 [11]. TEV_{ND}, PVY_{ND}, TuMV_{ND}: ND natural diversity. Simulations: four *in silico* replicates.

179 Regarding the simulated data, S and NS mutations are equally represented along potyvirus mutated
180 genomes, independently of the protein length ($R^2 \simeq 0.94$), thus validating our random model for a
181 number of mutations above 1000. As expected, it is indicative of the correlation between the protein
182 sequence length and the number of S or NS mutations obtained at random in the absence of any
183 biological bias.

184 To be noticed, for all simulated data, the number of NS mutations is 3 times higher than the number
185 of S mutations. By contrast, for the experimental data, the number of S and NS mutations is equivalent.
186 Assuming that there is little or no selection pressure on S mutations, we can extrapolate the number
187 of NS mutations before selection. So, in the TEV 2015 experiment [12], the total number of mutations
188 before selection would be 300 mutations distributed in 75 S and 225 NS mutations, 278 mutations for
189 the PVY 2015 experiment [10] and 1077 mutations for the PVY 2017 experiment [11]. We can see that
190 for two datasets, before selection, we do not reach the minimum required of 1000 mutations
191 previously defined to obtain a robust analysis.

192 The amount of S is rigorously proportional with the protein length irrespective of its function. This is
193 not the case for NS, and some proteins contain more mutations than others. For instance, it appears
194 that P1 significantly accumulates more NS than HC-Pro, P3, CI, N1b and CP ($P < 0.02$; Z test) (Figure 2
195 and Figure S2).

196
197

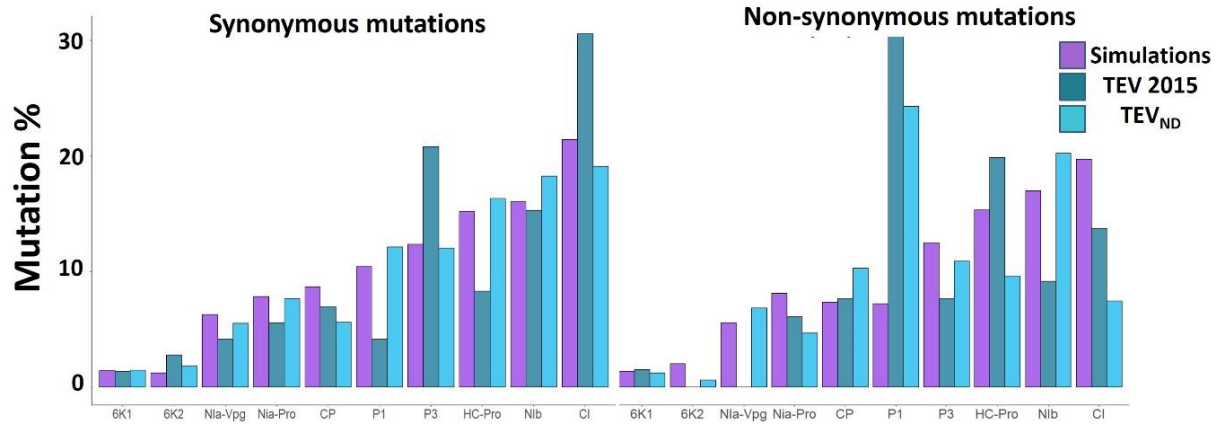


Figure 2. Mutation % in the TEV proteins from the experimental evolution [12], natural diversity and simulations. For PVY and TuMV see supplemental data. The proteins are sorted from the smallest to the largest, left to right: 6K1, 6K2, Nia-VPg, Nia-Pro, CP, P1, P3, Hc-Pro, Nib, CI.

198

199 DISTRIBUTION OF MUTATIONS NS and S in IDRs and ORs

200 In the second part of the study, the analysis was no longer conducted on individual proteins, but on
 201 all IDRs and ORs distributed along the coding sequences in the viral genomes. In order to analyze the
 202 distribution of each type of mutation in the IDRs or ORs, we defined the ratio R for synonymous
 203 mutations as:

204

$$R_s = \frac{\%IDR_S}{\%OR_S} \quad (1)$$

205 with $\%IDR_S$ and $\%OR_S$ defined as

206

$$\%IDR_S = \frac{\text{Number of } S \text{ mutations in IDR}}{\text{Total number mutations in IDR}} * 100 \quad (2)$$

$$\%OR_S = \frac{\text{Number of } S \text{ mutations in OR}}{\text{Total number mutations in OR}} * 100 \quad (3)$$

207

208 Equations 1-3 also apply for the calculation of R_{NS} , the ratio R for non-synonymous mutations.

209

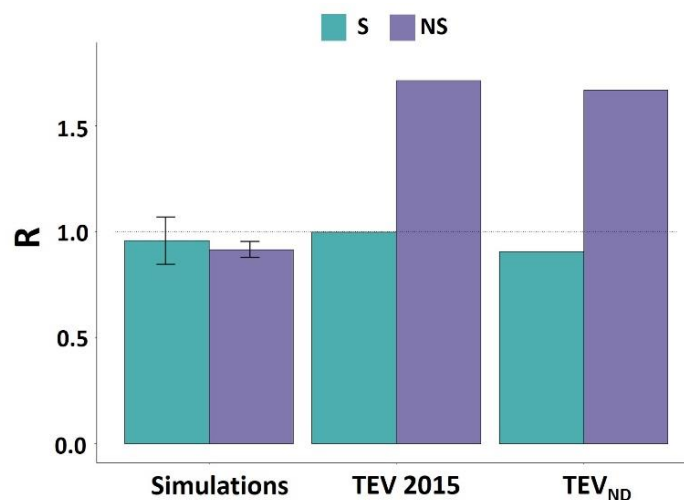


Figure 3. Ratio between the percentage of mutation (S or NS) present in IDRs and ORs for the TEV genome. For data sets from the two other studies [10,11], TuMV and PVY see supplemental data.

210 The ratio of synonymous mutations between IDRs and ORs deduced from the experimental data were
211 close to 1 (Figure 3 and Figure S3). This ratio is comparable to that obtained by simulation which mimics
212 random mutations and reflects the absence of impact of synonymous mutations at the protein level.
213 By contrast, for NS mutations, a large and significant difference between the experimental and
214 simulated data could be observed ($p < 0.02$, χ^2 for each of the four simulations, Table 2). Indeed, the
215 ratio higher than 1 observed in the case of the experimental data indicates an over-representation of
216 NS mutations within the IDRs compare to the ORs (Figure 3). In the case of PVY, this difference with
217 simulated mutations was only verified for the PVY 2017 dataset [11] (Table 2).

218

Simulations	TEV 2015	TEV _{ND}	PVY 2015	PVY 2017	PVY _{ND}	TuMV _{ND}
A	0,014	2.10-6	0,52	0,04	0,0005	3.10-8
B	0,026	1.10-5	0,46	0,03	0,0001	5.10-8
C	0,018	5.10-6	0,33	0,01	8.10-6	2.10-8
D	0,028	2.10-5	0,27	0,0049	1.10-7	6.10-9

Table 2. p values of χ^2 test for percentage of NS mutations in IDRs between simulated and experimental data (TEV 2015, PVY 2015, PVY 2017) or natural diversity (TEV_{ND}, PVY_{ND}, TuMV_{ND}). There was no experimental datasets available for TuMV. Significance, $p < 0.05$.

219

220 With respect to the analysis of natural diversity, no differences were observed between S mutations
221 within IDRs and ORs, in agreement with simulated data. By contrast, a significant over-representation
222 of NS was observed in the IDRs for all viruses (Figure 3 and Figure S3).

223

224 Altogether, those results indicate that IDRs are more prone to accumulate adaptive mutations than
225 more structured regions, at both short (experimental evolution) and longer (natural diversity)
226 evolutionary time-scales.

227

228 **Comparison of the physicochemical disturbance of amino acid substitutions in potyviral** 229 **intrinsically disordered versus ordered regions**

230 We also analyzed possible amino acid substitution biases between IDRs and ORs with respect to their
231 physico-chemical properties. To compare the physico-chemical nature of the substitutive amino acids
232 (NS mutations) in IDRs and ORs, we used the BLOSUM62 matrix [18]. This matrix uses the natural
233 diversity between very conserved regions of evolutionary divergent protein sequences. The set of
234 sequences is aligned to a reference sequence. For each position where a substitution occurs, the
235 probability of occurrence of each of the 19 other amino acids is calculated, resulting in score values
236 ranging between -4 and 11. The higher the score, the higher is the likelihood of substitution. It was
237 observed that the highest replacement probabilities is correlated to amino acids with similar physico-
238 chemical properties (charges, hydrophilicity-hydrophobicity, amino acid size) [19]. Upon amino acid
239 substitutions, the more drastic physicochemical changes are (the lower the score value), the more
240 destabilizing these changes are in terms of structure.

241 For each virus (TEV, PVY, and TuMV), we assessed whether the natural selection discriminated
242 differently within IDRs and ORs for amino acid substitutions with respect to their impact on biophysical
243 changes. For each type of region (IDRs or ORs), a comparative statistical analysis (Dunn test) was thus
244 performed between the natural diversity, experimental and simulated data sets (Table 2 and Table S4).
245 When amino acid substitutions occur in IDRs, their BLOSUM62 scores in the simulated data and those
246 in the PVY and TuMV data sets belong to the same statistical group. In the case of TEV, the natural
247 diversity data shows a slight difference with three of the four simulations. Importantly, its diversity
248 was represented by a set of only 6 genomes while those of PVY and TuMV were illustrated by 100
249 genomes each. In contrast, for all three potyviruses, regarding the amino acid substitutions present in
250 the ORs, both natural diversity and experimental evolution data have a significantly higher BLOSUM62

251 score than the simulated data. The high BLOSUM62 score observed as associated to ordered regions
 252 supports the idea that amino acids substitutions occurring in those regions are globally poorly
 253 destabilizing at the physicochemical and structural level. Reciprocally, with lower BLOSUM62 scores
 254 than the ones observed in ORs, IDRs would be more permissive to drastic physicochemical changes.
 255

A	Groups	
	IDRs	ORs
PVY 2015	abc	ab
PVY 2017	abc	a
PVY _{ND}	ac	a
Sim A	bc	c
Sim B	bc	c
Sim C	b	bc
Sim D	bc	c

B	Groups	
	IDRs	ORs
TEV 2015	ab	ab
TEV _{ND}	b	b
Sim A	a	a
Sim B	a	a
Sim C	a	a
Sim D	ab	a

C	Groups	
	IDRs	ORs
TuMV _{ND}	a	a
Sim A	a	b
Sim B	a	b
Sim C	a	b
Sim D	a	b

Table 2. Differences in physicochemical properties associated with amino acid substitutions were assessed using scores derived from the BLOSSOM62 substitution matrix. For each type of region (IDRs or ORs) groups (a,b and c) were determined by running a Dunn test (p value adjustment method: Bonferroni). For (A) PVY genome, (B) TEV genome and (C) TuMV genome.

256
 257 **Are NS mutations in IDRs driven toward the conservation of disorder promoting amino acids ?**
 258 We investigated whether the conservation of disorder during evolution could be a selection criterion
 259 using amino acid disorder promoting scores (see material and method section). Amino acid residues
 260 were grouped into order promoting, neutral or disorder promoting scores, ranging from 0 (amino acid
 261 most frequently present in ORs) to 1 (amino acid most frequently present in IDRs), were attributed to
 262 each of the 20 amino acids [17,20].

263 We first examined if, substitutions were preferentially targeting order or disorder promoting amino
 264 acids, and this, whether in IDRs or ORs. We did not observe any significant differences between
 265 biological data within either IDRs or ORs and simulated data (Table S5-A). We concluded that there is
 266 no natural tendency for evolution to target substitutions preferentially toward order or disorder
 267 promoting amino acids.

268 Then, we considered the possibility that non-synonymous mutations (NS) could preferably give order
 269 or disorder promoting amino acids (Table S5-B). We observed unbiased random substitution, and this,
 270 both in ORs or IDRs, in accordance with simulations. Finally, we aimed at assessing a possible tendency
 271 for substitution by amino acids that are more prone to promote order or disorder (Table S5-C). At each
 272 position where a NS mutation was observed, we calculated the difference in promoter score between
 273 the amino acid in the reference genome and the replacing amino acid in each of the genomes
 274 describing the diversity in the biological data. Again, we did not observe significant differences
 275 between naturally selected and simulated mutations. We did not detect any differences between
 276 biological and simulated data in the promoter score for the synonymous mutations, either in the IDRs
 277 or ORs. However, global disorder is generally conserved during evolution [21–23], and more specifically
 278 in RNA viruses [4,8]. Therefore, the analysis of local substitutions does not reflect this evolutionary
 279 trend that can be observed globally at the scale of a protein region. It turns out that the analysis of
 280 substitutions in terms of physico-chemical modulations sounds more relevant than the use of the
 281 order-disorder promoter scale.

282
 283
 284
 285

286 Discussion

287 *Mutational robustness differences between IDRs and ORs.*

288 Potyviruses constitute, together with begomoviruses, the two largest viral genera described to date
289 among plant viruses [24,25]. Potyviruses are very damaging to field crops and embrace a very wide
290 host range [26]. Most of them are generalists and as such, provide a rich model for studying viral
291 adaptation. In this study, we tested the hypothesis that among these viruses, mutational robustness
292 was greater in the disordered regions of their proteomes than in the ordered regions. We analyzed the
293 distribution of mutations in the genomes of potyviruses belonging to three different species, PVY, TEV
294 and TuMV, representative of the genus. The datasets used included viral genomes resulting from both
295 short evolutionary scale (experimental evolution) and longer evolutionary scales, with the use of
296 natural diversity. An analysis of the two datasets showed that IDRs and ORs are subject to different
297 evolutionary mechanisms, with disordered regions evolving towards significantly more amino acid
298 polymorphism than ordered regions. The selection pressure that applies to ordered regions thus tends
299 towards a conservative evolution while that which applies to disordered regions rather supports a
300 divergent evolution. It is quite easy to understand the evolutionary mechanism at work in the ordered
301 regions. These protein regions have a strong structure-function relationship. At the molecular level,
302 these regions are defined by geometries of constrained atomic interactions with few degrees of
303 freedom and well-packed hydrophobic cores. These regions have significantly higher BLOSUM62
304 scores than would result from random substitutions. Substituted amino acids have physicochemical
305 natures close to those of the original amino acids. Conversely, the lower topological requirement in
306 disordered regions results in substitutions close to the random substitution pattern. Within ordered
307 regions, mutations compensate each other to prevent instability according to an epistatic model. On
308 the long term, such compensation leads to changes in sequence and function (protein evolvability)
309 [6,27]. In these regions the selection pressure strongly operates to preserve function which results in
310 amino acids conservation. In disordered regions, the notion of structural stability is less relevant and
311 amino acids substitutions may have less functional impact [4]. This could constitute an alternative
312 model for protein evolvability, presumably on a shorter evolutive timeline consistent with the rapid
313 adaptation characteristic of viruses. From an evolutionary standpoint, the dogma of the structure-
314 function relationship (conservation of function requiring conserved structures and therefore close
315 substitutions) requires to be tempered in the case of IDRs.

316

317 *Does amino acid polymorphism in potyvirus proteome IDRs undergoes positive selection?*

318 We examined the hypothesis that intrinsic disorder could be selected to generate a pool of mutations
319 available for adaptive function. To obtain the diversity observed in IDRs, two successive processes,
320 namely the generation of mutations and their selection, are involved. The first process can be favored
321 by codon volatility. Codon volatility is defined as the proportion of a codon's point mutation neighbors
322 that code for different amino acids [28]. We investigated whether the nucleotide sequences encoding
323 IDRs used codons of higher volatility than the sequences of ORs, thus favoring the generation of non-
324 synonymous mutations. We did not observe greater codon volatility in the disordered regions than in
325 the ORs that could explain the greater amino acid polymorphism observed in the disordered regions.
326 Thus, with respect to the volatility criterion, we have no evidence to support that amino acid
327 polymorphism in disordered regions undergoes positive selection, generating a potential adaptive pool
328 for the virus. Although amino acid polymorphism in these regions may participate in potyvirus
329 adaptation, the conservation of intrinsic disorder during evolution is the result primarily of the second
330 process, a selection pressure dictated by the essential biochemical functions it performs to ensure
331 virus replication in the host. In any case, the mutational permissiveness and diversity that arise from
332 the selection of structure-function relationships within IDRs is likely to favor the adaptive potential of
333 the virus. It cannot therefore be excluded that IDRs are also selected according to this last criterion,
334 even if this hypothesis remains difficult to assess.

335

336 *Evolutionary features of nucleotide sequences encoding IDRs*

337 It should be expected that the synonymous codon usage pattern of viruses would be shaped by
338 selecting specific codon subsets to match the most abundant host transfer RNAs (tRNAs). However,
339 the codon usage of many viruses is very different from the optimal codons present in the host [29].
340 Interestingly, it was recently reported that codon usage in virus IDRs is less optimized for the host than
341 in ORs [30]. In the case of NS mutations, this is in line with our observation that IDRs are more robust
342 to mutations than OR, and thus evolve faster. This prevents fixation of codons optimized for the host
343 (figure 5). The preservation of codon diversity in these regions may also provide a reservoir for a faster
344 adaptation of the viruses to various hosts.
345

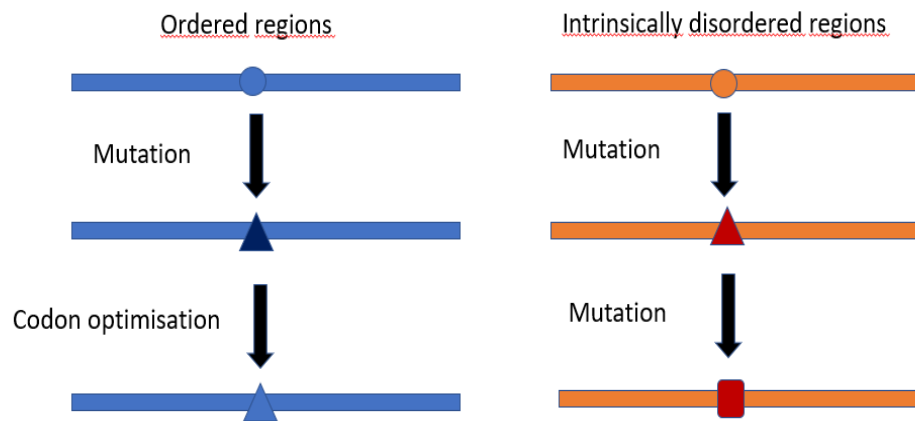


Figure 5. Mutational robustness of IDRs and codon optimization. The low rate of non-synonymous mutations in the ORs allows the optimization of the sequence towards the selection of abundant codons in the host. The high rate of non-synonymous mutations in IDRs prevents this optimization.

346
347 Because of the presence of less frequent codons in IDRs, the corresponding pools of loaded tRNAs in
348 the host cell are lower than those of abundant codons. Consequently, translational dynamics is likely
349 to be slowed down when the ribosome machinery enters a mRNA sequence encoding for disordered
350 regions [31,32]. This may result in an instability of the translation product [33]. IDRs are generally taken
351 over either co-translationally or post-translationally by chaperones. This handling does not favor the
352 selection of optimized codons and contributes to the preservation of amino acid polymorphism in IDRs.
353 There is an intricate interplay of molecular chaperones and protein disorder in the evolvability of
354 protein networks [34].

355 Taken all together, the data obtained unambiguously show that potyvirus IDRs and ODRs follow very
356 different evolutive paths with respect to their mutational robustness. These results force to consider
357 the hypothesis that during selection, adaptive solutions could emerge from the amino acid
358 polymorphism carried by IDRs.
359

360 **Bibliography**

- 361 1. Bellay J, Michaut M, Kim T, Han S, Colak R, Myers CL, et al. An omics perspective of protein
362 disorder. *Mol Biosyst.* 2012;8: 185–193. doi:10.1039/c1mb05235g
- 363 2. Habchi J, Tompa P, Longhi S, Uversky VN. Introducing Protein Intrinsic Disorder. *Chem Rev.*
364 2014;114: 6561–6588. doi:10.1021/cr400514h
- 365 3. Tokuriki N, Oldfield CJ, Uversky VN, Berezovsky IN, Tawfik DS. Do viral proteins possess unique
366 biophysical features? *Trends Biochem Sci.* 2009;34: 53–59. doi:10.1016/j.tibs.2008.10.009
- 367 4. Gitlin L, Hagai T, LaBarbera A, Solovey M, Andino R. Rapid Evolution of Virus Sequences in
368 Intrinsically Disordered Protein Regions. *PLoS Pathog.* 2014;10.

- 369 doi:10.1371/journal.ppat.1004529
- 370 5. Walter J, Charon J, Hu Y, Lachat J, Leger T, Lafforgue G, et al. Comparative analysis of
371 mutational robustness of the intrinsically disordered viral protein VPg and of its interactor
372 eIF4E. *PLoS One*. 2019;14: e0211725. doi:10.1371/journal.pone.0211725
- 373 6. Bershtein S, Segal M, Bekerman R, Tokuriki N, Tawfik DS. Robustness-epistasis link shapes the
374 fitness landscape of a randomly drifting protein. *Nature*. 2006;444: 929–932.
375 doi:10.1038/nature05385
- 376 7. Elena SF, Agudelo-Romero P, Carrasco P, Codoner FM, Martin S, Torres-Barcelo C, et al.
377 Experimental evolution of plant RNA viruses. *Heredity (Edinb)*. 2008;100: 478–483.
378 doi:10.1038/sj.hdy.6801088 ET - 2008/02/07
- 379 8. Charon J, Theil S, Nicaise V, Michon T. Protein intrinsic disorder within the Potyvirus genus:
380 From proteome-wide analysis to functional annotation. *Mol Biosyst*. 2016;12: 634–652.
381 doi:10.1039/c5mb00677e
- 382 9. Charon J, Barra A, Walter J, Millot P, Hébrard E, Moury B, et al. First Experimental Assessment
383 of Protein Intrinsic Disorder Involvement in an RNA Virus Natural Adaptive Process. *Mol Biol
384 Evol*. 2018;35: 38–49. doi:10.1093/molbev/msx249
- 385 10. Kutnjak D, Rupar M, Gutierrez-Aguirre I, Curk T, Kreuze JF, Ravnikar M. Deep Sequencing of
386 Virus-Derived Small Interfering RNAs and RNA from Viral Particles Shows Highly Similar
387 Mutational Landscapes of a Plant Virus Population. *J Virol*. 2015;89: 4760–4769.
388 doi:10.1128/jvi.03685-14
- 389 11. Kutnjak D, Elena SF, Ravnikar M. Time-Sampled Population Sequencing Reveals the Interplay
390 of Selection and Genetic Drift in Experimental Evolution of Potato Virus Y. *J Virol*. 2017;91.
391 doi:10.1128/jvi.00690-17
- 392 12. Cuevas JM, Willemsen A, Hillung J, Zwart MP, Elena SF. Temporal dynamics of intrahost
393 molecular evolution for a plant RNA virus. *Mol Biol Evol*. 2015;32: 1132–1147.
394 doi:10.1093/molbev/msv028
- 395 13. Schadt EE, Turner S, Kasarskis A. A Window into Third Generation Sequencing. *Hum Mol
396 Genet*. 2010;19: 227–240. doi:10.1093/hmg/ddq416
- 397 14. Sickmeier M, Hamilton JA, LeGall T, Vacic V, Cortese MS, Tantos A, et al. DisProt: the Database
398 of Disordered Proteins. *Nucleic Acids Res*. 2007;35: D786–D793. doi:10.1093/nar/gkl893
- 399 15. Xue B, Dunbrack RL, Williams RW, Dunker AK, Uversky VN. PONDR-FIT: A meta-predictor of
400 intrinsically disordered amino acids. *Biochim Biophys Acta - Proteins Proteomics*. 2010;1804:
401 996–1010. doi:10.1016/j.bbapap.2010.01.011
- 402 16. Tromas N, Elena SF. The Rate and Spectrum of Spontaneous Mutations in a Plant RNA Virus.
403 *Genetics*. 2010;185: 983–989. doi:10.1534/genetics.110.115915
- 404 17. Radivojac P, Iakoucheva LM, Oldfield CJ, Obradovic Z, Uversky VN, Dunker AK. Intrinsic
405 disorder and functional proteomics. *Biophys J*. 2006;12/13. 2007;92: 1439–1456.
406 doi:10.1529/biophysj.106.094045
- 407 18. Henikoff S, Henikoff JG. Amino acid substitution matrices from protein blocks. *Proc Natl Acad
408 Sci*. 1992;89: 10915–10919. doi:10.1073/pnas.89.22.10915
- 409 19. Rudnicki WR, Mroczek T, Cudek P. Amino acid properties conserved in molecular evolution.
410 *PLoS One*. 2014;9: e98983–e98983. doi:10.1371/journal.pone.0098983
- 411 20. Dunker AK, Obradovic Z. The protein trinity—linking function and disorder. *Nat Biotechnol*.
412 2001;19: 805–6 ST – The protein trinity—linking function. doi:10.1038/nbt0901-805
413 nbt0901-805 [pii] ET - 2001/09/05
- 414 21. Chen JW, Romero P, Uversky VN, Dunker AK. Conservation of intrinsic disorder in protein
415 domains and families: I. A database of conserved predicted disordered regions. *J Proteome
416 Res*. 2006;5: 879–887. Available:
417 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&
418 list_uids=16602695](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=16602695)
- 419 22. Ahrens J, Dos Santos HG, Siltberg-Liberles J. The Nuanced Interplay of Intrinsic Disorder and
420 Other Structural Properties Driving Protein Evolution. *Mol Biol Evol*. 2016;33: 2248–2256.
421 Available: <http://dx.doi.org/10.1093/molbev/msw092>
- 422 23. Ahrens JB, Nunez-Castilla J, Siltberg-Liberles J. Evolution of intrinsic disorder in eukaryotic

- 423 proteins. *Cell Mol Life Sci.* 2017;74: 3163–3174. doi:10.1007/s00018-017-2559-0
- 424 24. Nicaise V. Crop immunity against viruses: outcomes and future challenges. *Front Plant Sci.*
- 425 2014;5. doi:10.3389/fpls.2014.00660
- 426 25. Scholthof K-BG, Adkins S, Czosnek H, Palukaitis P, Jacquot E, Hohn T, et al. Top 10 plant viruses
- 427 in molecular plant pathology. *Mol Plant Pathol.* 2011;12: 938–954. doi:10.1111/j.1364-
- 428 3703.2011.00752.x
- 429 26. Moury B, Desbiez C. Host range evolution of potyviruses: A global phylogenetic analysis.
- 430 *Viruses.* 2020;12. doi:10.3390/v12010111
- 431 27. Tokuriki N, Tawfik DS. Protein dynamism and evolvability. *Science* (80-). 2009/04/11.
- 432 2009;324: 203–207. doi:10.1126/science.1169375
- 433 28. Zhang J. On the evolution of codon volatility. *Genetics.* 2005;169: 495–501.
- 434 doi:10.1534/genetics.104.034884
- 435 29. Jitobaom K, Phakaratsakul S, Sirihongthong T, Chotewutmontri S, Suriyaphol P, Suptawiwat O,
- 436 et al. Codon usage similarity between viral and some host genes suggests a codon-specific
- 437 translational regulation. *Heliyon.* 2020;6. doi:10.1016/j.heliyon.2020.e03915
- 438 30. Kumar N, Kaushik R, Tennakoon C, Uversky VN, Longhi S, Zhang KYJ, et al. Insights into the
- 439 evolutionary forces that shape the codon usage in the viral genome segments encoding
- 440 intrinsically disordered protein regions. *Brief Bioinform.* 2021;22. doi:10.1093/bib/bbab145
- 441 31. Yu CH, Dang Y, Zhou Z, Wu C, Zhao F, Sachs MS, et al. Codon Usage Influences the Local Rate
- 442 of Translation Elongation to Regulate Co-translational Protein Folding. *Mol Cell.* 2015;59: 744–
- 443 754. doi:10.1016/j.molcel.2015.07.018
- 444 32. Yang Q, Yu CH, Zhao F, Dang Y, Wu C, Xie P, et al. ERF1 mediates codon usage effects on
- 445 mRNA translation efficiency through premature termination at rare codons. *Nucleic Acids Res.*
- 446 2019. doi:10.1093/nar/gkz710
- 447 33. Mitarai N, Sneppen K, Pedersen S. Ribosome Collisions and Translation Efficiency:
- 448 Optimization by Codon Usage and mRNA Destabilization. *J Mol Biol.* 2008.
- 449 doi:10.1016/j.jmb.2008.06.068
- 450 34. Pechmann S, Frydman J. Interplay between chaperones and protein disorder promotes the
- 451 evolution of protein networks. *PLoS Comput Biol.* 2014;10: e1003674.
- 452 doi:10.1371/journal.pcbi.1003674
- 453