# Phospho-islands and the evolution of phosphorylated amino acids in mammals

**Mikhail Moldovan** [Corresp., 1] , **Mikhail S Gelfand** [1, 2]

[1] Skolkovo Institute of Science and Technology, Moscow, Russia

[2] A. A. Kharkevich Institute for Information Transmission Problems, Moscow, Russia

Corresponding Author: Mikhail Moldovan
Email address: mika.moldovan@gmail.com

## Background

Protein phosphorylation is the best studied post-translational modification strongly influencing protein function. Phosphorylated amino acids not only differ in physico-chemical properties from non-phosphorylated counterparts, but also exhibit different evolutionary patterns, tending to mutate to and originate from negatively charged amino acids. The distribution of phosphosites along protein sequences is non-uniform, as phosphosites tend to cluster, forming so-called phospho-islands.

## Methods

Here, we have developed an HMM-based procedure for the identification of phospho-islands and studied the properties of the obtained phosphorylation clusters. To check robustness of evolutionary analysis, we consider different models for the reconstructions of ancestral phosphorylation states.

## Results

Clustered phosphosites differ from individual phosphosites in several functional and evolutionary aspects including underrepresentation of phosphotyrosines, higher conservation, more frequent mutations to negatively charged amino acids. The spectrum of tissues, frequencies of specific phosphorylation contexts, and mutational patterns observed near clustered sites also are different.

1 **Phospho-islands and the evolution of phosphorylated**
2 **amino acids in mammals**
3

4 **Mikhail A. Moldovan[1], Mikhail S. Gelfand[1,2]**
5 [1] Skolkovo Institute of Science and Technology, Bolshoy Boulevard 30, bld. 1, Moscow, Russia
6 121205
7 [2] A.A.Kharkevich Institute for Information Transmission Problems (RAS), Bolshoy Karetny Per.
8 19, bld.1, Moscow, Russia 127051

9
10 Corresponding Author:
11 Mikhail A. Moldovan[1]
12 Bolshoy Boulevard 30, bld. 1, Moscow, 121205, Russia
13 Email address: mika.moldovan@gmail.com
14

# Abstract

**Background**

Protein phosphorylation is the best studied post-translational modification strongly influencing protein function. Phosphorylated amino acids not only differ in physico-chemical properties from non-phosphorylated counterparts, but also exhibit different evolutionary patterns, tending to mutate to and originate from negatively charged amino acids. The distribution of phosphosites along protein sequences is non-uniform, as phosphosites tend to cluster, forming so-called phospho-islands.

**Methods**

Here, we have developed an HMM-based procedure for the identification of phospho-islands and studied the properties of the obtained phosphorylation clusters. To check robustness of evolutionary analysis, we consider different models for the reconstructions of ancestral phosphorylation states.

**Results**

Clustered phosphosites differ from individual phosphosites in several functional and evolutionary aspects including underrepresentation of phosphotyrosines, higher conservation, more frequent mutations to negatively charged amino acids. The spectrum of tissues, frequencies of specific phosphorylation contexts, and mutational patterns observed near clustered sites also are different.

# Introduction

Protein post-translational modifications (PTMs) are important for a living cell (Schweiger and Linial 2010; Kurmangaliyev et al. 2011; Studer et al. 2016; Huang et al. 2017). By changing physico-chemical properties of proteins, PTMs affect their function, often introducing novel biological features (Pearlman et al. 2011). To date, hundreds of thousands of PTMs in various organisms have been identified and various databases containing information about PTMs have been compiled (Ptacek and Snyder, 2006; Huang et al. 2017).

Protein phosphorylation is likely both the most common and the best studied PTM (Ptacek and Snyder, 2006; Schweiger and Linial, 2010; Huang et al. 2017). Phosphorylation introduces a negative charge and a large chemical group to the local protein structure, hence strongly affecting the protein conformation (Pearlman et al. 2011; Nishi et al. 2014). As a result, diverse cellular signaling pathways are based on sequential phosphorylation events (Moses and Landry 2010; Pearlman et al. 2011; Ardito et al. 2017). In eukaryotes, phosphorylation sites (phosphosites) are mainly represented by serines, threonines, and tyrosines (which we here refer to as STY amonoacids), with only a minor fraction involving other amino acids, such as histidine (Fuhs and Hunter 2017; Huang et al. 2017).

Phosphosites are overrepresented in disordered regions (DRs) of proteins, i.e. in regions devoid of secondary or tertiary structure, usually located on the surface of a protein globule (Iakoucheva 2004). Hence, studies of the evolution of phosphosites have mainly concentrated

58    on sites located in DRs (Kurmangaliyev et al. 2011; Miao et al. 2018). In particular, it has been
59    shown, that phosphosites tend to arise from negatively charged amino acids (NCAs) more
60    frequently than their non-phosphorylated counterparts, and, in a number of cases, retain
61    structural features initially maintained by NCAs (Kurmangaliyev et al. 2011; Miao et al. 2018).
62    As phosphorylation is often highly conserved (Macek et al. 2007), experimental limitations on
63    the number of model species with established phosphosites may be overcome in evolutionary
64    studies by formally assigning phosphorylation labels to homologous sites (Kurmangaliyev et al.
65    2011; Huang et al. 2017). However, this approach requires a degree of caution when dealing
66    with evolutionary trees of substantial depths, *e.g.* only a small fraction of yeast phosphosites are
67    conserved between species separated by ~1400 My (million years), while about a half of
68    phosphosites are conserved at a shorter time (~360 My) (Studer et al. 2016). At smaller
69    distances, this method may be applied to infer some evolutionary properties of phosphosites,
70    *e. g.* in *Drosophila* or in vertebrate species, phosphosites tend to mutate to NCA
71    (Kurmangaliyev et al. 2011; Miao et al. 2018).
72           Phosphorylation can be both a constitutive modification and a way to transiently modify
73    the protein function (Landry et al. 2014). In the former case, the change of a phosphosite to
74    NCA should not cause a significant fitness reduction, as physico-chemical properties are not
75    strongly affected, whereas in the latter case a mutation would have dire consequences (Moses
76    and Landry 2010; Landry et al. 2014).
77           In proteins, phosphosites often form co-localized groups called phosphorylation islands
78    or phosphorylation clusters, and about a half of phosphorylated serines and threonines are
79    located in such clusters (Schweiger and Linial 2010). While individual phosphosites function as
80    simple switches, phospho-islands are phosphorylated in a cooperative manner, so that the
81    probability of a phosphorylation event at a focal site strongly depends on the phosphorylation of
82    adjacent sites, and when the number of phosphorylated amino acids exceeds a threshold, the
83    cumulative negative charge of the phosphate groups introduces functionally significant changes
84    to the protein structure (Landry et al. 2014).
85           Accurate procedures for the identification of phosphosites and next-generation
86    sequencing technologies yielded large numbers of well-annotated phosphosites (Altenhoff et al.
87    2017; Huang et al. 2017; UniProt 2018) enabling us to develop an accurate automatic
88    procedure for the identification of phospho-site clusters we call phospho-islands. We show that
89    clustered phosphosites exhibit evolutionary properties distinct from those of individual
90    phosphosites, in particular, an enhanced mutation rate to NCA and altered mutational patterns
91    of amino acids in the phosphosite vicinity. Our study complements earlier observations on the
92    general evolutionary patterns in phosphosites with the analysis of mutations in non-serine
93    phosphosites and the demonstration of differences in the evolution of clustered and individual
94    phosphorylated residues.
95

## Materials & Methods
97    **Data**
98           The phosphosite data for human, mouse and rat proteomes were downloaded from the
99    iPTMnet database (Huang et al. 2017). The phosphorylation breadth values for the mouse
100   dataset were obtained from (Huttlin et al. 2010). Human, mouse and rat proteomes were
101   obtained from the UniProt database (UniProt 2018). Vertebrate orthologous gene groups

102 (OGGs) for human and mouse proteomes were downloaded from the OMA database (Altenhoff
103 et al. 2017). Then, all paralogous sequences and all non-mammalian sequences were excluded
104 from the obtained OGGs.
105

106 **Alignments and Trees**
107       We searched for homologous proteins in three proteomes with pairwise BLASTp
108 alignments (Altschul et al. 1990). Pairs of proteins with highest scores were considered closest
109 homologs. The information about closest homologs was subsequently used to predict HMR
110 phosphosites. OGG were aligned by the ClustalO multiple protein alignment (Sievers et al.
111 2014) and, while the HMR phosphosites were identified based on Muscle pairwise protein
112 alignments (Edgar 2004). The mammalian phylogenetic tree was obtained from Timetree
113 (Kumar et al. 2017).
114

115 **Phosphorylation retention upon mutations**
116       After the identification of homologous protein pairs in human/mouse and mouse/rat
117 proteomes and the proteome alignment construction, we identified homologous phosphosites as
118 homologous STY resudues which were shown to be phosphorylated in both species. We have
119 shown that phosphorylation is retained on S-T and T-S mutation by comparing two pairs of
120 ratention probabilities (Fig. 1C): $p(pS\text{-}pS)$ with $p(pS\text{-}pT \mid S)$ and $p(pT\text{-}pT)$ with $p(pS\text{-}pT \mid T)$
121 (analogously for the phosphorylation of tyrosines), $p(pX\text{-}pX)$ being defined as the fraction of X
122 amino acids phosphorylated in both considered species:
123

124
$$p(pX - pX) = \frac{\#(pX - pX)}{\#(pX - pX) + \#(pX - X)}$$

125

126 and $p(pX_1\text{-}pX_2)$, as the fraction of phosphorylated $X_1$ residues in one species given that in
127 another species another amino acid residue ($X_2$) is also phosphorylated:
128

129
$$p(pX_1 - pX_2 | X_1) = \frac{\#(pX_1 - pX_2)}{\#(pX_1 - pX_2) + \#(pX_1 - X_2)}$$

130

131
$$p(pX_1 - pX_2 | X_2) = \frac{\#(pX_1 - pX_2)}{\#(pX_1 - pX_2) + \#(X_1 - pX_2)}$$

132

133 Homologous phosphosite lists from the human/mouse and human/rat pairs were merged to
134 produce HMR phosphosite list of human phosphosites.
135

136 **False-positive rates of phosphorylation identification by homologous propagation**
137       We assessed the quality of the phosphorylation prediction via homologous propagation
138 approaches by counting false-positive rates of phosphosite predictions in species with large
139 phosphosite lists. As the numbers of predicted phosphosites drastically differed between
140 species (Huang et al. 2017), we considered multiway predictions in each case as characteristics
141 of the procedure performance. Hence, considering mouse phosphosite predicted by homology
142 with known human phosphosites, we also considered human phosphosites predicted based on

143 known mouse phosphosites. The false-positive rate was assessed as the proportion of correctly
144 predicted phosphosites among the STY aminoacids in one species homologous to phosphosite
145 positions in other considered species.
146       When assessing the quality of phosphosite predictions based on phosphosites
147 experimentally identified in at least two species, we considered human, mouse. and rat and the
148 lists of phosphosites homologous between human and mouse and between human and rat. In
149 these cases, predictions were made for rat and mouse, respectively with the false-positive rate
150 assessed by the same approach as in the previous case.
151

## Mutation matrices

153       To obtain single-aminoacid mutation matrices, we first reconstructed ancestral states
154 with the PAML software (Yang 2007). For the reconstruction, we used OGG alignments which
155 did not contain paralogs and pruned mammalian trees retaining only organisms contributing to
156 corresponding OGG alignments. The alignment of both extant and reconstructed ancestral
157 sequences and the corresponding trees were then used to construct mutation matrices, where
158 we distinguished the phosphorylated and non-phosphorylated states of STY amino acids. Here,
159 the phosphorylation state was assigned to STY amino acids using the phosphorylation
160 propagation approach described above. When calculating the mutation matrix, we did not count
161 mutations predicted to happen on branches leading from the root to first-order nodes, as PAML
162 did not reconstruct them well without an outgroup (Koshi and Goldshtein 1996; Yang 2007).
163 Tree pruning and calculating the mutation matrix count were implemented in *ad hoc* python
164 scripts using functions from the ete3 python module.
165

## Disordered regions and phospho-island prediction

167       Disordered regions were predicted with the PONDR VSL2 software with default
168 parameters (Xue et al. 2010). Phosho-islands were predicted with the Viterbi algorithm (Viterbi
169 1967). Emission probabilities for the algorithm were obtained as the ratio of density values in the
170 $S$ distribution decomposition (likelihood ratio normalized to 1) (Fig. 2a). Transitional probabilities
171 were set to 0.2 to maximize the likeness of obtained distribution of $S$ within phospho-islands and
172 the one predicted by the decomposition procedure (Fig. 2b).
173

## Phosphosite contexts

175       We employed the list of phosphosite contexts as well as the binary decision-tree
176 procedure to define the context of a given phosphosite from Villen et al. 2007. The procedure is
177 as follows. (i) Proline context is assigned if there is a proline at position +1 relative to the
178 phosphosite. (ii) Acidic context is assigned if there are five or six E/D amino acids at positions
179 +1 to +6 relative to the phosphosite. (iii) Basic context is assigned if there is a R/K amino acid at
180 position –3. (iv) Acidic context is assigned if there are D/E amino acids at any of positions +1,
181 +2 or +3. (v) Basic context is assigned if there are at least two R/K amino acids at positions –6
182 to –1. Otherwise, no context is assigned and we denote this as the "O" (other) context. We
183 consider tyrosine phosphosites separately and formally assign the with the "Y" (tyrosine)
184 context.
185

## Local mutation matrices

187       We computed local substitution matrices (LSM) as the substitution matrices for amino
188 acids located within a frame with the radius $k$ centered at a phosphorylated serine or threonine.
189 When computing LSMs, we did not count mutations of or resulting in STY amino acids to
190 exclude the effects introduced by the presence and abundance of phospho-islands. We have
191 set $k$ to 1, 3, 5, and 7 and selected 5 as for this value we observed the strongest effect, that is,
192 obtained the largest number of mutations with frequencies statistically different from those for
193 non-phosphorylated serines and threonines.
194
195 **Statistics**
196       When comparing frequencies, we used the $\chi^2$ test if all values in the contingency matrix
197 exceeded 20 and Fisher's exact test otherwise. To correct for multiple testing, we used the
198 Bonferroni correction with the scaling factor set to 17 for the substitution vector comparison and
199 to 17×17 for the comparison of substitution matrices with excluded STY amino acids. 95% two-
200 tailed confidence intervals shown in figures were computed by the $\chi^2$ or Fisher's exact test. The
201 significance of obtained Pearson's correlation coefficients was assessed with the F-statistic.
202
203 **Code availability**
204       *Ad hoc* scripts were written in Python. Graphs were built using R. All scripts and data
205 analysis protocols are available online at https://github.com/mikemoldovan/phosphosites.
206
207 # Results
208 **Conserved phosphosites**
209       As protein phosphorylation in a vast majority of organisms has not been studied or has
210 been studied rather poorly (Huang et al. 2017), the evolutionary analyses of phosphosites
211 typically rely on the assumption of absolute conservation of the phosphorylation label assigned
212 to STY amino acids on a considered tree (Kurmangaliyev et al. 2011; Miao et al. 2018). Thus, if,
213 for instance, a serine is phosphorylated in human, we, following this approach, would consider
214 any mutation in the homologous position of the type S-to-X to be a mutation of a phosphorylated
215 serine to amino acid X (Fig. 1b). However, the comprehensive analysis of yeast phosphosites
216 has shown low conservation of the phosphorylation label at the timescales of the order 100 My
217 and more (Studer et al. 2016). Thus, we have considered only orthologous groups of
218 mammalian proteins, present in the OMA database (Altenhoff et al. 2017). The mammalian
219 phylogenetic tree is about 177 My deep (Kumar et al. 2017), which corresponds to about 50% of
220 the phosphorylation loss in the 182 My-deep yeast *Saccharomyces-Lachancea* evolutionary
221 path (Studer et al. 2016). The tree contains three organisms with well-studied
222 phosphoproteomes: human (227834 sites), mouse (92943 sites), and rat (24466 sites) (Huang
223 et al. 2017) (Fig. 1a).
224       Still, the expected 50% of mispredicted phosphosites could render an accurate
225 evolutionary analysis impossible. This could be partially offset by considering phosphosites
226 conserved in well-studied lineages. Thus, we compiled a set of human phosphosites
227 homologous to residues phosphorylated also in mouse and/or rat, which we will further refer to
228 as human-mouse/rat (HMR) phosphosites. The HMR set consists of 53437 sites covering
229 54.6% and 61.2% of known mouse and rat phosphosites, respectively, which is consistent with

230    the above-mentioned observation about 50% phosphorylation loss in yeast on evolutionary
231    distances similar to the ones between the human and rodent lineages (Fig. 1ab).
232          We consider the HMR set to be enriched in accurately predicted phosphosites. Indeed,
233    by considering conserved phosphosites, we substantially reduce the number of mispredictions.
234    If we simply propagated human phosphorylation labels to mouse and *vice versa* we would get
235    about 77.6% and 42.3% of false positive labels, respectively. However, sites conserved
236    between human and rat or sites conserved between rat and mouse would yield about twofold
237    lesser percentages of 41.9% and 19.9% of false positives in mouse and human, respectively.
238    The obtained percentages can be considered as upper estimates of false positive rates, as
239    current experimental phosphosite coverage in mammals cannot guarantee the identification of
240    all conserved phosphosites (Huang et al. 2017). Thus, the HMR dataset is sufficiently robust for
241    the prediction of phosphorylation labels in less-studied mammalian lineages.
242          Treatment of STY amino acids homologous to phosphorylated ones as phosphorylated
243    yields another possible caveat, stemming from the possible loss of phosphorylation upon STY-
244    to-STY mutations. To assess this effect, we compared the probabilities of phosphosite retention
245    upon pSTY-to-STY mutation, pSTY indicating the phosphorylated state, and the respective
246    probabilities in the situation when a mutation has not occurred for a pair of species with well-
247    established phosphosite lists, i.e. human and mouse (Fig. 1c). We have observed only a minor,
248    insignificant decrease of the probabilities of the phosphorylation retention in the cases of pS-pT
249    and pS-pY mismatches relative to the pT-pT states in mouse and human, indicating the general
250    conservation of the phosphorylation label upon amino acid substitution. An interesting
251    observation here is that the pS-pS states appear to be the most conserved ones (Fig. 1c).
252    Taken together, these results indicate the evolutionary stability of phosphorylation states upon
253    mutation.
254          The increased evolutionary robustness of the pS state relative to the pT and pY states
255    should manifest as overrepresentation of phosphoserines among phosphosites with respect to
256    non-phosphorylated amino acid positions. Thus, we assessed the relative abundances of pSTY
257    amino acids in the HMR dataset relative to the established human phosphosite set and to the
258    set of non-phosphorylated STY amino acids. Serines and threonines, comprising the vast
259    majority of the pSTY amino acids, are, respectively, over- and underrepresented in the
260    phosphosite sets (Fig. 1cd). This effect is significantly more pronounced in the HMR dataset
261    relative to the total human phosphosite dataset, further supporting the observation about lower
262    conservation of pT relative to pS, as the HMR dataset is enriched in conserved phosphosites by
263    design.
264
265    **Phosphorylation islands**
266          The distribution of distances between phosphosites is different from that of randomly
267    chosen serines and threonines even accounting for the tendency of phosphosites to occur in
268    disordered regions (DRs) (Schweiger and Linial 2010) (Fig. 2a). However, this observation
269    depends on an arbitrary definition of phosphorylation islands as groups of phosphosites
270    separated by at most four amino acids (Schweiger and Linial 2010). We have developed an
271    approach that reduces the degree of arbitrariness in the definition of phospho-islands.
272          Let $S$ be the distribution of amino acid distances between adjacent phosphosites in DRs.
273    The logarithm of $S$ is not unimodal (Fig. 2a), and we suggest that it is a superposition of two

274    distributions: one generated by phosphosites in phospho-islands and the other reflecting
275    phosphosites outside phospho-islands (left and right peaks, respectively). The latter distribution
276    can be obtained from random sampling from DRs of non-phosphorylated STY amino acids while
277    preserving the amino acid composition and the sample size, as we expect individual
278    phosphosites to emerge independently while maintaining the preference towards DRs (Fig. 2c).
279    Gamma distribution has a good continuous fit to log($S$+1) for randomly sampled STY amino
280    acids located in DRs. Given its universality and low number of parameters (Friedman et al.
281    2006; Reiss et al. 2007; Mendoza-Parra et al. 2013), we have selected gamma distribution as a
282    reasonable model for log($S$+1) (Fig. 2c). Assuming that the distribution of log($S$+1) values for
283    phosphosites located in phospho-islands should belong to the same family and fixing the
284    parameters of the previously obtained distribution, we decomposed the distribution of log($S$+1)
285    values into the weighted sum of two gamma distributions, one of which corresponds to STYs
286    located in phospho-islands and the other one, to remaining STYs in DRs (Fig 2a, red and grey
287    curves, respectively). From these two gamma distributions we obtained parameters for a hidden
288    Markov model, which, in turn, was used to map phosphorylation islands. The distributions of $S$
289    values for phosphosites in identified islands and the distribution for other phosphosites yielded a
290    good match to the expected ones (Fig. 2bd).
291        Both for the HMR and mouse datasets, more than half of phosphosites are located in
292    phospho-islands (61% and 56%, respectively) (Fig. 2e, Suppl. Fig. S8AB). For human
293    phosphosites, however, we see a larger proportion of sites (53%) located outside phospho-
294    islands. In the latter case the distributions in the decomposition differ less, compared to the
295    former two cases (Fig. 2a, Suppl. Fig. S8). It could be caused by a larger density of
296    phosphosites in DRs of the human proteome, resulting from higher experimental coverage; that
297    would lead to generally lower $S$ values, which, in turn, could cause the right peak in the log($S$+1)
298    distribution to merge with the left peak, rendering the underlying gamma-distributions less
299    distinguishable. To validate this explanation, we randomly sampled 40% of human
300    phosphosites, so that the sample size matched the one for mouse phosphosites; however, the
301    results on this rarefied dataset did not change (Fig. 2e, Suppl. Fig. S8C) indicating that our
302    procedure is robust with respect to phosphosite sample sizes. Hence, phospho-islands for the
303    human dataset are identified with a lower accuracy than those for the HMR and mouse
304    datasets. This could be caused by different experimental technique applied to the human
305    phosphosites, compared to the one used for mouse and rat phosphosites, and by a possibly
306    large number of false-positive phosphosites in the former case (Huttlin et al. 2010; Bekker-
307    Jensen et al. 2017; Xu et al. 2017) (see Discussion).
308        In phospho-islands, the overall pSTY-amino acid composition differs from that of
309    individual phosphosites, mainly because the fraction of threonines is significantly higher in
310    phospho-islands at the expense of the lower fraction of tyrosines (Fig. 2f). Also, the
311    conservation of residues in phospho-islands is larger than that of the individual sites (Fig. 2g).
312    Overall, the general properties of clustered phosphosites seem to differ from those of individual
313    phosphosites.
314        We do not observe phospho-islands in ordered regions, as the distribution of log($S$+1)
315    values in this case seemingly cannot be decomposed into a weighted sum of two unimodal
316    distributions and is largely skewed to the left even relative to the distribution of log($S$+1) values
317    in phospho-islands. Hence, either virtually all pairs of adjacent sites there comprise dense

318    phospho-islands, or a more complex model possibly incorporating tertiary protein structure
319    features is required to infer phospho-islands in this case (Suppl. Fig. S8E).
320
321

322    **Mutational patterns of phosphorylated amino acids**
323        Next, we have reconstructed the ancestral states for all mammalian orthologous protein
324    groups not containing paralogs and calculated the proportions of mutations $P(X_1 \rightarrow X_2)$, where $X_1$
325    and $X_2$ are different amino acids. We treated phosphorylated and non-phosphorylated states of
326    STY amino acids as distinct states. We then introduced a measure of difference in mutation
327    rates for phosphorylated STY and their non-phosphorylated counterparts. For a mutation of an
328    STY amino acid $X$ to a non-STY amino acid $Z$ we define $R(X, Z) = P(pX \rightarrow Z) / P(X \rightarrow Z)$. If $X^*$ is
329    another STY amino acid, $R(X, X^*) = P(pX \rightarrow pX^*) / P(X \rightarrow X^*)$. Thus, the $R$ value for a given type
330    of mutations is the proportion of the considered mutation of a phosphorylated STY amino acid
331    among other mutations normalized by the fraction of respective mutations of the non-
332    phosphorylated STY counterpart. The $R$ values are thus not affected by differences in the
333    mutation rates between phosphorylated and non-phosphorylated amino acids, as all
334    probabilities are implicitly normalized by the mutation rates of $pX$ and $X$.
335        We firstly consider phosphosites located in DRs. For phosphorserines from the HMR
336    dataset we confirm earlier observations: phosphoserines mutate to NCA more frequently than
337    non-phosphorylated serines (Fig. 3a). The $R$ values for serine mutation to aspartate, $R(S,D)$,
338    and glutamate, $R(S,E)$, are both significantly larger than 1 (1.2, $p<0.01$ and 1.7, $p<0.001$,
339    respectively; $\chi^2$ test) and, interestingly, they differ substantially ($p<0.001$, multiple random
340    Poisson sampling test). Similarly, asparagine and glutamine $R$ values differ, with $R(S,N)=0.9$
341    ($p<0.001$, $\chi^2$ test), significantly lower than 1, and $R(S,Q)=1.4$ ($p<0.001$, $\chi^2$ test), significantly
342    higher than 1. The rate of mutation to lysine significantly differs for phosphorylated and non-
343    phosphorylated serines ($p<0.001$, $\chi^2$ test). Interestingly, the mutation rate to another positively
344    charged amino acid, arginine, is significantly lower than expected ($p<0.01$, $\chi^2$ test). For non-
345    polar amino acids generally no significant differences in the $R$ values between phosphorylated
346    and non-phosporylated serines are observed, but for methionine and proline, the calculated
347    values are significant: $R(S,M)>1$ ($p<0.001$, $\chi^2$ test) and $R(S,P)<1$ ($p<0.01$, $\chi^2$ test).
348        In earlier studies, only mutations of serines or to serines had been considered, as the
349    available data did not allow for statistically significant results for threonine and tyrosine
350    (Kurmangaliyev et al. 2011; Miao et al. 2018). Here, we see that phosphorylated threonines
351    from the HMR dataset tend to mutate to serines (Fig. 3b). At that, phosphorylated serines
352    mutate to threonines more frequently than their non-phosphorylated counterparts for all
353    considered samples, i.e. for the human, mouse and HMR sets (Fig. 3b, Suppl. Figs. S2-S7).
354    Phosphorylated tyrosines tend to avoid mutations to isoleucine ($p<0.05$, $\chi^2$ test) and, for human
355    samples, to arginine ($p<0.05$, $\chi^2$ test) and glycine ($p<0.001$, $\chi^2$ test) (Fig. 3b, Suppl. Figs. S2-
356    S7). Phospho-tyrosines in the mouse dataset show a weaker tendency for the avoidance of the
357    mutations to aspartate than the non-phosphorylated ones ($p<0.05$, $\chi^2$ test) while the rate of pY-
358    to-I mutations is higher (Fig. 3b).
359        Separate analysis of mutations in phospho-islands and in individual phosphosites yields
360    three observations. Firstly, alterations of mutation patterns of phosphoserines and
361    phosphothreonines (pST) in DRs relative to non-phosphorylated ST in DRs are similar to the

362     patterns observed for the clustered pST and, to a lesser extent, to those observed for individual
363     pSTs (Fig. 3b). This is mostly due to the fact that the mutational patterns of clustered pSTs
364     generally differ from those of their non-phosphorylated counterparts to a greater extent than the
365     mutational patterns of individual phosphoserines do (Fig. 3b, Suppl. Fig. S1). Secondly, for
366     phosphotyrosines, alterations in their mutational patterns brought about by phosphorylation are
367     mostly explained by individual phosphotyrosines. The mutational patterns of individual sites
368     deviate from the ones observed for non-phosphorylated tyrosines more than those of clustered
369     phosphotyrosines (Fig. 3b, Suppl. Figs. S2-S7). Also, if we compare the $R$ values calculated for
370     all possible mutations in clustered *vs.* individual phosphosites, the $R$ value corresponding to the
371     S-to-E mutation will be significantly higher for the set of clustered phosphosites ($p=0.009$, $\chi^2$
372     test, Suppl. Fig. S1). Hence, we posit that the general phosphosite mutational pattern alterations
373     can be explained mostly by mutations in clustered phosphosites for phosphoserines and
374     phosphothreonines and by individual sites when phosphotyrosines are considered.
375         We also studied mutation patterns in ordered regions (ORs), and observed that
376     phosphothreonines located in ORs demonstrate higher T-to-S mutation rates (Fig. 3b) relative to
377     those of non-phosphorylated threonines located in ORs. Also, sites located in ORs demonstrate
378     enhanced S-to-T and Y-to-T mutation rates relative to non-phosphorylated serines and
379     threonines in ORs, respectively (Fig. 3b).
380

381     **Phosphosite contexts**
382         Sequence contexts of phosphosites generally fall into three categories: acidic (A), basic
383     (B), and proline (P) motifs, with tyrosine phosphosites comprising a special class (Y) (Villen et
384     al. 2007; Huttlin et al. 2010). For each phosphosite from each dataset we have identified its
385     context. As in previous studies (Villen et al. 2007; Huttlin et al. 2010), phosphosites not
386     assigned with any of these context classes were considered as having "other" (O) motif. We
387     studied the distribution of these motifs for all classes of phosphosites.
388         In DRs, relative to ORs, we observed a higher percentage of phosphosites with assigned
389     contexts (Fig. 4a). P-phosphosites demonstrate the highest overrepresentation in DRs, with
390     25% of DR phosphosites having the proline motif. Phospho-islands, compared to individual
391     phosphosites, contain more phosphosites with assigned motifs relative to individual
392     phosphosites. In DRs, there are more B- and P-phosphosites and fewer A-phosphosites and Y-
393     phosphosites among clustered sites than among individual ones.
394

395     **Phosphorylation breadth**
396         An important feature of a phosphosite is the "phosphorylation breadth", that is, the
397     number of tissues where it is phosphorylated. In this study, the maximal phosphorylation
398     breadth is nine, as the phosphorylation data for nine mouse tissues are available (Huttlin et al.
399     2010). Among broadly expressed phosphosites (present in all nine tissues), compared to tissue-
400     specific ones (present in only one tissue), very few sites have unassigned contexts (O) and
401     almost none are tyrosine phosphosites. The fraction of acidic phosphosites (24%) is
402     substantially lower among tissue-specific sites relative to broadly phosphorylated ones (37%)
403     ($p<0.001$, $\chi^2$ test) (Fig. 4a).
404         As mentioned above, the pS-to-E mutation yields the highest value, $R(S,E)$ (Fig. 3a) and
405     represents the only mutation with significantly different $R$ values in phospho-islands and

406    individual sites ($p$=0.009, χ² test, Suppl. Fig. S1). At that, $R$(S,E) significantly increase with
407    increasing breadth of expression (Fig. 4b), from $R$(S,E)=1.14 for tissue-specific phosphosites to
408    $R$(S,E)=6.64 for broadly expressed phosphosites ($p$=0.016, t-test).
409            Finally, we compared percentages of phosphosites with different breadths in ORs *vs.*
410    DRs and in phospho-islands *vs.* individual phosphosites (Fig. 4cd). As the phosphorylation
411    breadth increases, so does the fraction of clustered phosphosites, reaching 85% for sites
412    phosphorylated in nine tissues; the fraction of phosphosites in DRs also increases, reaching
413    95.4%.
414            Hence, broadly expressed phosphosites have well-defined motifs, tend towards
415    disordered regions and to phospho-islands, have mostly acidic context, and mutate to NCA
416    more frequently than tissue-specific phosphosites.
417
418    **Mutation patterns in the proximity of phosphosites**
419            We now show that not only phoshosites require special motifs (Huttlin et al. 2010), but
420    the mutational context of clustered phosphosites differs from that of individual sites. To assess
421    evolutionary dynamics associated with phosphosite motifs, we analyzed mutational patterns in
422    ±3 amino acid windows of HMR ST phosphosites located in DRs and compared them with those
423    of non-phosphorylated ST amino acids. The ±3 window was selected, as it yielded the strongest
424    effect in terms of the number of mutations with rates statistically distinct from the expected ones
425    (Suppl. Fig. S9AB). We did not consider phosphotyrosines, as they have not been shown to
426    possess any discernible general motif apart from the phosphorylated tyrosine itself (Huttlin et al.
427    2010).
428            We introduce the measure $Q$ defined as $Q\left(X_1^p \rightarrow X_2^p\right) = P(X_1^p \rightarrow X_2^p)/P(X_1^n \rightarrow X_2^n)$, where $X_1^p$
429    and $X_2^p$ are amino acids near phosphorylated serines and threonines and $X_1^n$ and $X_2^n$ are amino
430    acids near non-phosphorylated serines and threonines. $Q$ measures overrepresentation of a
431    given mutation in the proximity of pST amino acids relative to ST amino acids. We also
432    considered sites located in phospho-islands and individual phosphosites separately (Fig. 5,
433    Suppl. Fig. S9CD).
434            In the whole HMR dataset, 22 types of non-phosphorylated amino acid substitutions out
435    of the total of 289 have $Q$ values statistically different from the expected value 1 ($p$<0.05, χ² test
436    with the Bonferroni correction), among them three pairs of mutually reverse mutations (Fig. 5).
437    As expected from the conservation of the phosphosite contexts, mutations between positively
438    and negatively charged amino acids, potentially changing acidic to basic contexts and *vice*
439    *versa*, are underrepresented, whereas E-to-D, D-to-E and K-to-R, not changing the context
440    type, are overrepresented. The P-to-A substitution is overrepresented, thus indicating the
441    instability of proline contexts. Interestingly, all three mutations with $Q$ values exceeding 2.5
442    involve lysine, two of them being reverse mutations F-to-K and K-to-F. The fourth most
443    overrepresented mutation, Y-to-G with $Q$(Y→G)=2.5, could explain the lack of tyrosine
444    phosphosites in DRs, as a large fraction of DR phosphosites are clustered with the distances
445    between sites not exceeding three amino acids. Thus, a large $Q$(Y→G) value would lead to
446    general underrepresentation of tyrosines in DRs.
447    Types of mutations with significant $Q$ values generally differ near clustered and individual
448    phosphosites (Suppl. Fig. S9CD). E-to-D, not changing the local acidic context type (Huttlin et

449    al. 2010), is overrepresented and E-to-K, disrupting the acidic context (Huttlin et al. 2010), is

450    underrepresented in both cases. On the other hand, around individual phosphosites, $Q(F{\rightarrow}K)=3.4$

451    and $Q(P{\rightarrow}A)=1.12$, indicating an enhanced birth rate of the basic context and disruption of the

452    proline context, respectively. The R-to-D mutation, disrupting the local basic context, also is

453    overrepresented near individual phosphosites. In general, among seven overrepresented

454    mutations near clustered phosphosites, only the K-to-P mutation disrupts the local basic context

455    in favor of the proline context and among seven overrepresented mutations near individual

456    phosphosites, three mutations (E-to-F, R-to-D, and P-to-A) could be regarded as context-

457    disrupting. Hence, the individual phosphosite contexts are somewhat less evolutionary stable and

458    thus the lower percentage of individual phosphosites with identifiable contexts might be due to

459    specific local context-disrupting mutation patterns for these phosphosites.

460

## Discussion

461

**Clustered vs. individual phosphosites**

462

463         We have demonstrated that clustered phosphosites differ from non-clustered ones in a

464    number of aspects: (i) overrepresentation of phosphothreonines and underrepresentation of

465    phosphotyrosines in phospho-islands (Fig. 2f); (ii) stronger conservation of clustered

466    phosphoserines and phosphothreonines (Fig. 2g); (iii) larger proportion of sites phosphorylated

467    in many tissues (Fig. 4C); (iv) significantly larger probability of mutations to glutamate for

468    clustered relative to the individual phosphoserines; (v) larger fraction of sites with specific motifs

469    in phospho-islands (Fig. 4A); (vi) mutational patterns in the proximity of phosphosites consistent

470    with the context-retention hypothesis (Fig. 5). What are possible explanations for the observed

471    effects?

472         Underrepresentation of phosphotyrosines in phospho-islands could be explained by

473    phosphorylation of clustered phosphosites being co-operative. As serines and threonines are

474    more similar to each other in their tendency to being phosphorylated by similar enzymes than

475    they are to tyrosine (Villen et al. 2007; Huttlin et al. 2010; Landry et al. 2014; Studer et al. 2016),

476    one would expect phospho-tyrosines to disrupt co-operative phosphorylation of adjacent ST

477    amino acids by being phosphorylated independently, thus introducing a negative charge which

478    would affect phosphorylation probabilities of the neighbouring amino acids (Landry et al. 2014).

479    Hence phospho-tyrosines could have been purged by selection from pST clusters.

480         Secondly, phosphosites located in phospho-islands are more conserved than individual

481    ones (Fig. 2g), as opposed to an earlier hypothesis that individual phosphosites are more

482    conserved than their clustered counterparts (Landry et al. 2014). Our result seems to contradict

483    to the notion that the cellular function of phosphosites in an island depends on the number of

484    phosphorylated residues rather than specific phosphorylated sites, whereas individual

485    phosphosites operate as single-site switches and hence should be more conserved (Landry et

486    al. 2014). However, this argument implies that phosphorylation of most individual phosphosites

487    is important for the organism's fitness, which may be not true (Landry et al. 2014; Miao et al.

488    2018) and hence our results do not contradict the model of evolution of functionally important

489    phosphosites.

490         Overrepresentation of phosphosites with defined motifs among the clustered ones (Fig.

491    4A) and reduced numbers of mutations disrupting the local contexts of the clustered sites

492   (Suppl. Fig. S9CD) may indicate enhanced selective pressure on clustered phoshosites and
493   their contexts. An indirect support of this claim comes from the overrepresentation of
494   ubiquitously phosphorylated sites among the clustered ones (Fig. 4c). Indeed, broad
495   phosphorylation requires a stronger local context and indicates the reduced probability of a
496   phosphosite being detected simply due to the noise inherent to the phosphorylation machinery
497   (Landry et al. 2014).
498         Mutations of phosphoserines located in DRs to NCA are generally overrepresented
499   among all mutations of the type pS-to-X relative to the corresponding mutations of non-
500   phosphorylated serines (Fig. 3b). This effect is stronger for clustered phosphosites and for
501   ubiquitously phosphorylated sites. Together with the observation  about clustered phosphosites
502   being on average more broadly phosphorylated than the individual ones, this suggests that a
503   large fraction of phosphosite clusters might be phosphorylated (nearly) constitutively, and thus
504   changes of individual phospho-serines to NCAs could experience lesser degrees of negative
505   selection acting upon the corresponding mutations, as these mutations introduce smaller
506   degrees of local electric charge shifts on the protein globule than the mutations of non-
507   phosphorylated serines to NCAs do.
508
509   **Two types of mutations**
510         In all considered phosphosite datasets, we have observed two types of pSTY-to-X
511   mutations overrepresented relative to STY-to-X mutations (Fig. 3b): (i) pSTY-to-pSTY,
512   especially pT-to-pS mutation and (ii) pSTY-to-NCA, especially pS-to-E mutations. The former
513   effect could be explained by the relaxed selection against pST-to-pST mutations due to the
514   phosphorylation machinery often not distinguishing between serines and threonines (Huttlin et
515   al. 2010; Miao et al. 2018). The overrepresentation of pT-to-pS mutation for all datasets,
516   including sites located in ORs, could stem from the higher probability of phosphosite retention
517   following a pT-to-pS mutation relative to the probability of phosphorylated threonine retention
518   when no mutations have occurred (Fig. 1c). Thus, the observed enhanced pT-to-pS mutation
519   rate could be due to the enhanced evolutionary stability of serine phosphorylation relative to the
520   threonine phosphorylation.
521         The enhanced serine-to-NCA mutation rates could stem from the physico-chemical
522   similarity of phosphorylated serines and negatively charged amino acids: both types of residues
523   introduce negatively charged groups of similar size to the protein globule. Thus, if
524   phosphorylation is (almost) constitutive, i.e. happens very frequently in a large number of
525   tissues, we would expect the serine-to-NCA mutation rate to be enhanced. Indeed, ubiquitous
526   phosphorylated serines have the pS-to-E mutation rate more than six-fold larger than the S-to-E
527   mutation rate (Fig. 4b). However, the same pattern does not hold for phospho-threonines (Fig.
528   3B).
529
530   **Human phosphosites**
531         The results obtained for the human set of phosphosites differ somewhat from those for
532   the mouse and HMR sets, like in cases with different STY amino acids representation among
533   phosphorylated amino acids (Fig. 1D), proportion of phosphosites located in phospho-islands
534   (Fig. 2E) or some mutational patterns of phosphorylated STY amino acids (Fig. 3B). This could
535   be explained by differences in experimental procedures used to obtain phosphosite lists for

536   human and for mouse and rat. Whereas for classic laboratory organisms, phosphosites are
537   obtained directly from the analysis of an organism or an analysis of its live organ (Huttlin et al.
538   2010), for human phosphosite inference immortalized cell lines, such as HeLa, are used
539   (Bekker-Jensen et al. 2017; Xu et al. 2017), with conditions differing from those *in vivo*, and
540   hence one could expect different patterns of phosphorylation. In particular, the lower rate of
541   mutations to NCA could be explained by overrepresentation of sites with noisy phosphorylation
542   manifesting only in cell lines under the conditions of experiments. The mutation of such a
543   residue to NCA would most likely result in the deleterious effect of an average non-
544   phosphorylated serine mutation to NCA (Jin and Pawson 2012). Thus, we propose that
545   phosphosites conserved between human and rodent lineages, called here HMR sites, are more
546   robust with respect to experimental techniques, and hence are better suited for phosphosite
547   evolutionary studies.
548
549   **Evolution of non-studied phosphosite groups**
550         Previous studies dedicated to the evolution of phosphosites have focused on
551   phosphoserines located in DRs. The large datasets employed in the present study enabled us
552   to assess the patterns of phosphothreonines, phosphotyrosines and sites located in ORs. Apart
553   from the largely enhanced pT-to-pS mutation proportions relative to T-to-S ones (Fig. 3b) no
554   patterns with straightforward biological explanation were observed in these cases. However, an
555   interesting observation here is the consistent, significantly enhanced rate of pY-to-I mutations
556   relative to the Y-to-I mutations in the mouse and HMR datasets (Fig. 3b).
557
558    **Perspectives**
559   We propose a simple yet accurate homology-based approach for the ancestral phosphosite
560   inference yielding in our case the set of HMR phosphosites. As the predicted fractions of
561   phosphorylation labels falsely assigned to internal tree nodes are much smaller than the ones for
562   other phosphosite datasets, HMR set poses a valuable source of data for evolutionary studies. A
563   practical extension of our homology-based approach could be a phosphosite prediction procedure
564   incorporating additional pieces of information such as the tendency of phosphosites to cluster,
565   the local phosphosite contexts, and the tree structure into the probabilistic model, which would
566   predict phosphosites with a high degree of accuracy. On the other hand, it would be interesting to
567   infer the interplay between phosphorylation and selection using population-genetics data.
568
569

# References

571   2018. UniProt: a worldwide hub of protein knowledge. Nucleic Acids Research [Internet]
572   47:D506–D515. Available from: http://dx.doi.org/10.1093/nar/gky1049
573   Altenhoff AM, Glover NM, Train C-M, Kaleb K, Warwick Vesztrocy A, Dylus D, de Farias TM,
574   Zile K, Stevenson C, Long J, et al. 2017. The OMA orthology database in 2018: retrieving
575   evolutionary relationships among all domains of life through richer web and programmatic
576   interfaces. Nucleic Acids Research [Internet] 46:D477–D485. Available from:
577   http://dx.doi.org/10.1093/nar/gkx1019

578 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.
579 Journal of Molecular Biology [Internet] 215:403–410. Available from:
580 http://dx.doi.org/10.1016/S0022-2836(05)80360-2
581 Ardito F, Giuliani M, Perrone D, Troiano G, Muzio LL. 2017. The crucial role of protein
582 phosphorylation in cell signaling and its use as targeted therapy (Review). International Journal
583 of Molecular Medicine [Internet] 40:271–280. Available from:
584 http://dx.doi.org/10.3892/ijmm.2017.3036
585 Bekker-Jensen DB, Kelstrup CD, Batth TS, Larsen SC, Haldrup C, Bramsen JB, Sørensen KD,
586 Høyer S, Ørntoft TF, Andersen CL, et al. 2017. An Optimized Shotgun Strategy for the Rapid
587 Generation of Comprehensive Human Proteomes. Cell Systems [Internet] 4:587–599.e4.
588 Available from: http://dx.doi.org/10.1016/j.cels.2017.05.009
589 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
590 throughput. Nucleic Acids Research [Internet] 32:1792–1797. Available from:
591 http://dx.doi.org/10.1093/nar/gkh340
592 Friedman, N., Cai, L., & Xie, X. S. (2006). Linking Stochastic Dynamics to Population
593 Distribution: An Analytical Framework of Gene Expression. Physical Review Letters, 97(16).
594 https://doi.org/10.1103/physrevlett.97.168302
595 Fuhs, S. R., & Hunter, T. (2017). pHisphorylation: the emergence of histidine phosphorylation as
596 a reversible regulatory modification. Current Opinion in Cell Biology, 45, 8–16.
597 https://doi.org/10.1016/j.ceb.2016.12.010
598 http://etetoolkit.org/
599 Huang H, Arighi CN, Ross KE, Ren J, Li G, Chen S-C, Wang Q, Cowart J, Vijay-Shanker K,
600 Wu CH. 2017. iPTMnet: an integrated resource for protein post-translational modification
601 network discovery. Nucleic Acids Research [Internet] 46:D542–D550. Available from:
602 http://dx.doi.org/10.1093/nar/gkx1104
603 Huttlin EL, Jedrychowski MP, Elias JE, Goswami T, Rad R, Beausoleil SA, Villén J, Haas W,
604 Sowa ME, Gygi SP. 2010. A Tissue-Specific Atlas of Mouse Protein Phosphorylation and
605 Expression. Cell [Internet] 143:1174–1189. Available from:
606 http://dx.doi.org/10.1016/j.cell.2010.12.001
607 Iakoucheva LM. 2004. The importance of intrinsic disorder for protein phosphorylation. Nucleic
608 Acids Research [Internet] 32:1037–1049. Available from: http://dx.doi.org/10.1093/nar/gkh253
609 Jin J, Pawson T. 2012. Modular evolution of phosphorylation-based signalling systems.
610 Philosophical Transactions of the Royal Society B: Biological Sciences [Internet] 367:2540–
611 2555. Available from: http://dx.doi.org/10.1098/rstb.2012.0106
612 Koshi JM, Goldstein RA. 1996. Probabilistic reconstruction of ancestral protein sequences.
613 Journal of Molecular Evolution [Internet] 42:313–320. Available from:
614 http://dx.doi.org/10.1007/BF02198858
615 Kumar S, Stecher G, Suleski M, Hedges SB. 2017. TimeTree: A Resource for Timelines,
616 Timetrees, and Divergence Times. Molecular Biology and Evolution [Internet] 34:1812–1819.
617 Available from: http://dx.doi.org/10.1093/molbev/msx116

618  Kurmangaliyev YZ, Goland A, Gelfand MS. 2011. Evolutionary patterns of phosphorylated
619  serines. Biology Direct [Internet] 6:8. Available from: http://dx.doi.org/10.1186/1745-6150-6-8
620  Landry CR, Freschi L, Zarin T, Moses AM. 2014. Turnover of protein phosphorylation evolving
621  under stabilizing selection. Frontiers in Genetics [Internet] 5. Available from:
622  http://dx.doi.org/10.3389/fgene.2014.00245
623  Macek B, Gnad F, Soufi B, Kumar C, Olsen JV, Mijakovic I, Mann M. 2007. Phosphoproteome
624  Analysis ofE. coliReveals Evolutionary Conservation of Bacterial Ser/Thr/Tyr Phosphorylation.
625  Molecular & Cellular Proteomics [Internet] 7:299–307. Available from:
626  http://dx.doi.org/10.1074/mcp.M700311-MCP200
627  Mendoza-Parra, M.-A., Nowicka, M., Van Gool, W., & Gronemeyer, H. (2013). Characterising
628  ChIP-seq binding patterns by model-based peak shape deconvolution. BMC Genomics, 14(1),
629  834. https://doi.org/10.1186/1471-2164-14-834
630  Miao B, Xiao Q, Chen W, Li Y, Wang Z. 2018. Evaluation of functionality for serine and
631  threonine phosphorylation with different evolutionary ages in human and mouse. BMC
632  Genomics [Internet] 19. Available from: http://dx.doi.org/10.1186/s12864-018-4661-6
633  Moses AM, Landry CR. 2010. Moving from transcriptional to phospho-evolution: generalizing
634  regulatory evolution? Trends in Genetics [Internet] 26:462–467. Available from:
635  http://dx.doi.org/10.1016/J.Tig.2010.08.002
636  Nishi, H., Shaytan, A., & Panchenko, A. R. (2014). Physicochemical mechanisms of protein
637  regulation by phosphorylation. Frontiers in Genetics, 5.
638  https://doi.org/10.3389/fgene.2014.00270
639  Pearlman SM, Serber Z, Ferrell JE Jr. 2011. A Mechanism for the Evolution of Phosphorylation
640  Sites. Cell [Internet] 147:934–946. Available from: http://dx.doi.org/10.1016/j.cell.2011.08.052
641  PTACEK J, SNYDER M. 2006. Charging it up: global analysis of protein phosphorylation.
642  Trends in Genetics [Internet] 22:545–554. Available from:
643  http://dx.doi.org/10.1016/j.tig.2006.08.005
644  Reiss, D. J., Facciotti, M. T., & Baliga, N. S. (2007). Model-based deconvolution of genome-
645  wide DNA binding. Bioinformatics, 24(3), 396–403.
646  https://doi.org/10.1093/bioinformatics/btm592
647  Schweiger R, Linial M. 2010. Cooperativity within proximal phosphorylation sites is revealed
648  from large-scale proteomics data. Biology Direct [Internet] 5:6. Available from:
649  http://dx.doi.org/10.1186/1745-6150-5-6
650  Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert
651  M, Soding J, et al. 2014. Fast, scalable generation of high-quality protein multiple sequence
652  alignments using Clustal Omega. Molecular Systems Biology [Internet] 7:539–539. Available
653  from: http://dx.doi.org/10.1038/msb.2011.75
654  Studer RA, Rodriguez-Mias RA, Haas KM, Hsu JI, Vieitez C, Sole C, Swaney DL, Stanford LB,
655  Liachko I, Bottcher R, et al. 2016. Evolution of protein phosphorylation across 18 fungal species.
656  Science [Internet] 354:229–232. Available from: http://dx.doi.org/10.1126/science.aaf2144

657 Villen J, Beausoleil SA, Gerber SA, Gygi SP. 2007. Large-scale phosphorylation analysis of
658 mouse liver. Proceedings of the National Academy of Sciences [Internet] 104:1488–1493.
659 Available from: http://dx.doi.org/10.1073/pnas.0609836104
660 Viterbi A. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding
661 algorithm. IEEE Transactions on Information Theory [Internet] 13:260–269. Available from:
662 http://dx.doi.org/10.1109/TIT.1967.1054010
663 Xu H, Chen X, Ying N, Wang M, Xu X, Shi R, Hua Y. 2017. Mass spectrometry-based
664 quantification of the cellular response to ultraviolet radiation in HeLa cells.Huen MS-Y, editor.
665 PLOS ONE [Internet] 12:e0186806. Available from:
666 http://dx.doi.org/10.1371/journal.pone.0186806
667 Xue B, Dunbrack RL, Williams RW, Dunker AK, Uversky VN. 2010. PONDR-FIT: A meta-
668 predictor of intrinsically disordered amino acids. Biochimica et Biophysica Acta (BBA) -
669 Proteins and Proteomics [Internet] 1804:996–1010. Available from:
670 http://dx.doi.org/10.1016/j.bbapap.2010.01.011
671 Yang Z. 2007. PAML 4: Phylogenetic Analysis by Maximum Likelihood. Molecular Biology
672 and Evolution [Internet] 24:1586–1591. Available from:
673 http://dx.doi.org/10.1093/molbev/msm088
674 Yu Y, Zhou H, Kong Y, Pan B, Chen L, Wang H, Hao P, Li X. 2016. The Landscape of A-to-I
675 RNA Editome Is Shaped by Both Positive and Purifying Selection.Schierup MH, editor. PLOS
676 Genetics [Internet] 12:e1006191. Available from:
677 http://dx.doi.org/10.1371/journal.pgen.1006191
678

679 **Figure captions**

680

681 **Figure 1 | Phosphosites considered in the study. (A)**. Venn diagram of iPTMnet human,
682 mouse and rat phosphosites. Intersections correspond to conserved phosphosites. The HMR
683 phosphosite dataset is shown in pink. **(B)**. Phosphosite assignment procedures. Given a tree of
684 a mammalian orthologous gene group and a column in the respective alignment, we assign
685 phosphorylation labels to ancestral and extant amino acids, firstly, by propagating labels from
686 one species to all other species in the tree (shown as separate red and blue arrows) and,
687 secondly, by propagating labels predicted both in the selected species (e.g. human, as shown)
688 and in one of the remaining species (mouse and rat); this corresponds to blue and red arrows
689 entering a given node in the tree. Phosphosites obtained by the latter procedure are referred to
690 as the HMR phosphosite dataset. In both procedures, phosphorylation is considered to be
691 retained both for direct and indirect STY-to-STY mutations. **(C)**. Retention of phosphorylation
692 upon mutation. Bars represent the probability of a conserved modification for the human dataset
693 in the case of mutation and if mutation has not occurred. The letter after the vertical bar is an
694 amino acid over which the probability was normalized. Three asterisks represent $p<0.001$ ($\chi^2$
695 test). **(D)**. STY amino acid content of three groups of phosphosite datasets.

696

697

698 **Figure 2 | Phospho-islands for the HMR phosphosite dataset. (A)**. The distribution of $log_{10}(S$
699 $+ 1)$ values (pink histogram) and its decomposition in two gamma distributions: the one for
700 phospho-islands (red curve) and for individual phosphosites (red curve). **(B)**. The distribution of
701 $log_{10}(S + 1)$ values for phosphosites predicted to be in phospho-islands. **(C)**. $log_{10}(S + 1)$ values
702 for non-phosphorylated STY amino acids randomly sampled from DRs with the same sample
703 size and amino acid content as in the HMR dataset. **(D)**. $log_{10}(S + 1)$ values for predicted
704 individual phosphosites. **(E)**. Numbers of individual phosphosites and sites in phospho-islands
705 for four datasets. **(F)**. Amino acid content of phospho-islands and individual phosphosites. **(G)**.
706 Frequency of mutations for phosphosites and individual amino acids. Asterisks depict
707 significantly different values ($p < 0.001$, $\chi^2$ test).

708
709
710 **Figure 3 | pX$_0$ -> X$_1$ substitution vectors. (A)**. $R$ values of the pS$\rightarrow X$ substitutions for serines
711 from the HMR dataset located in DRs. **(B)**. Substitution probabilities for phosphorylated STY
712 amino acids significantly different from those for non-phosphorylated STY amino acids for
713 several datasets. The significance levels are shown with the colors introduced in the panel in
714 (A). Abbreviations on the horizontal axis: ISL – phosphosites located in phospho-islands, IND –
715 individual phosphosites. DR – phosphosites from disordered regions, OR – phosphosites from
716 ordered regions.

717
718
719 **Figure 4 | Phosphosite contexts and phosphorylation breadth. (A)**. Overrepresentation of
720 phosphosite contexts in ordered *vs*. disoredred regions, in phospho-islands *vs*. individual
721 phosphosites and for broadly *vs*. narrowly distributed phosphorylated amino acids. One asterisk
722 and three asterisks indicate statistical significance at the levels of 0.05 and 0.001 respectively
723 ($\chi^2$ test). **(B)**. The dependence of $R$(pS, E) on the phosphosite breadth. Pearson's $r^2$ is equal to
724 0.53 with the t-test $p=0.016$. **(C)**. The dependence of phosphosite fraction in phospho-islands on
725 the phosphorylation breadth ($p=9*10^{-41}$, $\chi^2$ test). **(D)**. Percent of phosphosites in disordered
726 regions *vs*. phosphosite breadth ($p=4.1*10^{-10}$, $\chi^2$ test).

727
728
729 **Figure 5 | *Q* values of mutations near ST phosphosites with probabilities significantly**
730 **different from the expected ones.** Solid red lines connect mutually reverse mutations. Dashed
731 lines indicate quazy-reverse mutations of amino acids with common chemical properties.
732

# Figure 1
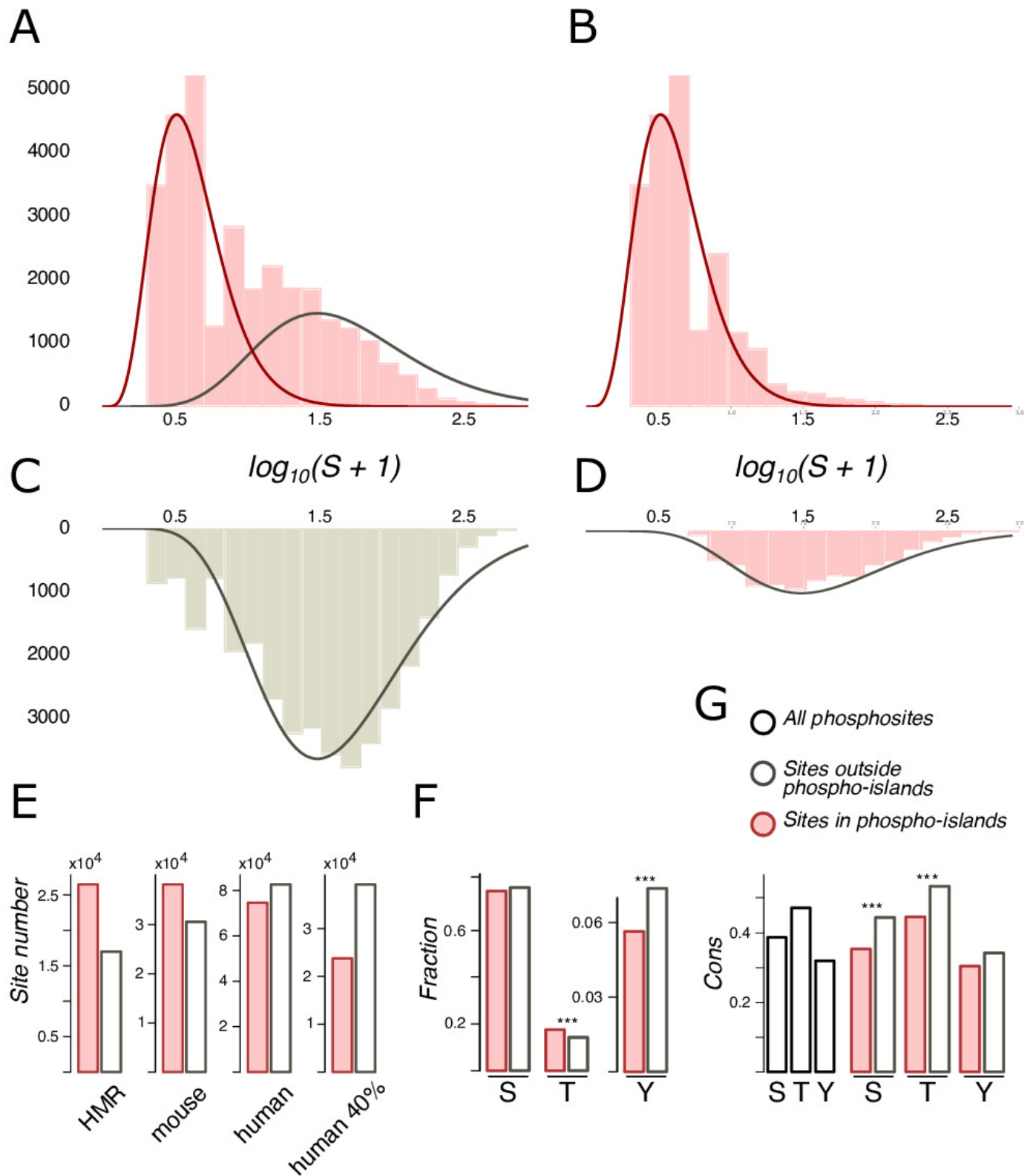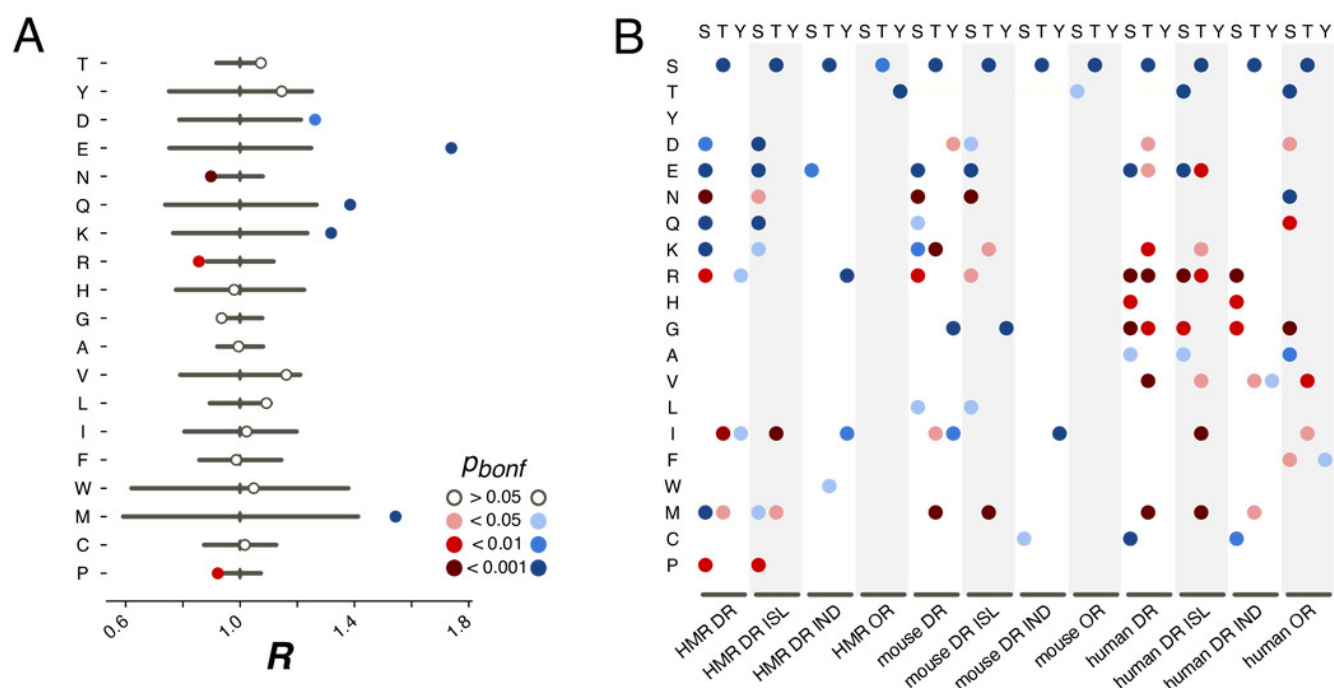
Figure 1 | Phosphosites considered in the study.

**(A)**. Venn diagram of iPTMnet human, mouse and rat phosphosites. Intersections correspond to conserved phosphosites. The HMR phosphosite dataset is shown in pink. **(B)**. Phosphosite assignment procedures. Given a tree of a mammalian orthologous gene group and a column in the respective alignment, we assign phosphorylation labels to ancestral and extant amino acids, firstly, by propagating labels from one species to all other species in the tree (shown as separate red and blue arrows) and, secondly, by propagating labels predicted both in the selected species (e.g. human, as shown) and in one of the remaining species (mouse and rat); this corresponds to blue and red arrows entering a given node in the tree. Phosphosites obtained by the latter procedure are referred to as the HMR phosphosite dataset. In both procedures, phosphorylation is considered to be retained both for direct and indirect STY-to-STY mutations. **(C)**. Retention of phosphorylation upon mutation. Bars represent the probability of a conserved modification for the human dataset in the case of mutation and if mutation has not occurred. The letter after the vertical bar is an amino acid over which the probability was normalized. Three asterisks represent $p<0.001$ ($\chi^2$ test). **(D)**. STY amino acid content of three groups of phosphosite datasets.

# Figure 2

Figure 2 | Phospho-islands for the HMR phosphosite dataset.

**(A)**. The distribution of $log_{10}(S + 1)$ values (pink histogram) and its decomposition in two gamma distributions: the one for phospho-islands (red curve) and for individual phosphosites (red curve). **(B)**. The distribution of $log_{10}(S + 1)$ values for phosphosites predicted to be in phospho-islands. **(C)**. $log_{10}(S + 1)$ values for non-phosphorylated STY amino acids randomly sampled from DRs with the same sample size and amino acid content as in the HMR dataset. **(D)**. $log_{10}(S + 1)$ values for predicted individual phosphosites. **(E)**. Numbers of individual phosphosites and sites in phospho-islands for four datasets. **(F)**. Amino acid content of phospho-islands and individual phosphosites. **(G)**. Frequency of mutations for phosphosites and individual amino acids. Asterisks depict significantly different values ($p < 0.001$, $\chi^2$ test).

# Figure 3

Figure 3 | pX$_0$ -> X$_1$ substitution vectors.

**(A)**. *R* values of the pS☐X substitutions for serines from the HMR dataset located in DRs. **(B)**. Substitution probabilities for phosphorylated STY amino acids significantly different from those for non-phosphorylated STY amino acids for several datasets. The significance levels are shown with the colors introduced in the panel in (A). Abbreviations on the horizontal axis: ISL – phosphosites located in phospho-islands, IND – individual phosphosites. DR – phosphosites from disordered regions, OR – phosphosites from ordered regions.
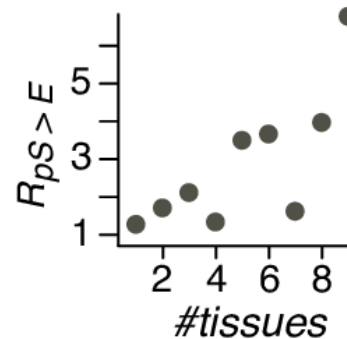
# Figure 4

Figure 4 | Phosphosite contexts and phosphorylation breadth.

**(A)**. Overrepresentation of phosphosite contexts in ordered *vs.* disoredred regions, in phospho-islands *vs.* individual phosphosites and for broadly *vs.* narrowly distributed phosphorylated amino acids. One asterisk and three asterisks indicate statistical significance at the levels of 0.05 and 0.001 respectively ($\chi^2$ test). **(B)**. The dependence of $R(pS,E)$ on the phosphosite breadth. Pearson's $r^2$ is equal to 0.53 with the t-test *p*=0.016. **(C)**. The dependence of phosphosite fraction in phospho-islands on the phosphorylation breadth ($p=9*10^{-41}$, $\chi^2$ test). **(D)**. Percent of phosphosites in disordered regions *vs.* phosphosite breadth ($p=4.1*10^{-10}$, $\chi^2$ test).
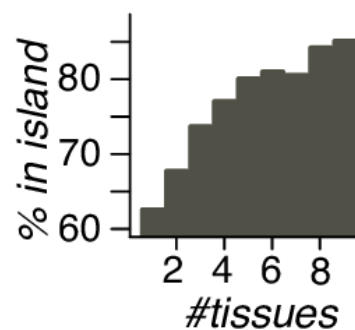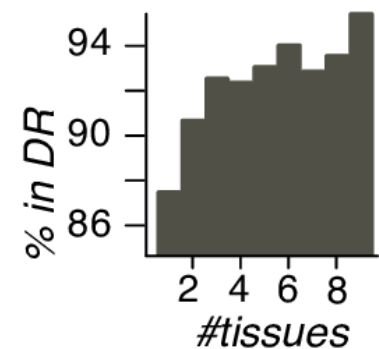
## A

# Figure 5

Figure 5 | *Q* values of mutations near ST phosphosites with probabilities significantly different from the expected ones.

Solid red lines connect mutually reverse mutations. Dashed lines indicate quazy-reverse mutations of amino acids with common chemical properties.