

DISCOVERY NOTES

Open Access

Evolutionary patterns of phosphorylated serines

Yerbol Z Kurmangaliyev^{1,2}, Alexander Goland¹, Mikhail S Gelfand^{1,3*}

Abstract: Posttranslationally modified amino acids are chemically distinct types of amino acids and in terms of evolution they might behave differently from their non-modified counterparts. In order to check this possibility, we reconstructed the evolutionary history of phosphorylated serines in several groups of organisms. Comparisons of substitution vectors have revealed some significant differences in the evolution of modified and corresponding non-modified amino acids. In particular, phosphoserines are more frequently substituted to aspartate and glutamate, compared to non-phosphorylated serines.

Reviewers: This article was reviewed by Arcady Mushegian and Sandor Pongor.

Findings

Post-translational modifications play an important role in diversifying protein structure and function [1,2]. Protein phosphorylation is one of the most important and widely distributed types of post-translational modifications. In eukaryotes, reversible protein phosphorylation plays a key role in the signal transduction and other processes [3,4]. Recent advances in mass spectrometry allowed for large-scale identification of phosphorylation events [5]. Analyses of these data have already revealed some specific structural and evolutionary features of phosphoserines. Phosphoserines tend to occur in intrinsically disordered regions [6-8] and regions corresponding to alternatively spliced gene segments [9]. Phosphorylated amino acids are more conserved than their non-phosphorylated counterparts [7,10-12]. Some very old phosphorylation events potentially can be common to organisms from *Archaea* to human [10].

Here we investigated another evolutionary aspect of protein modification sites. Since modified amino acids chemically are a distinct type of amino acids, in terms of evolution they might behave differently from their non-modified counterparts (on the top of the different level of conservation). To analyse differences in the evolution of standard amino acids and their modified counterparts, we reconstructed the evolution of phosphorylated amino acids in three groups of organisms. Particularly, we studied phosphorylation of serine in the human, fruit fly and yeast proteomes.

Phosphorylation sites were downloaded from the PHOSIDA [7] and PhosphoPEP [13] databases. For yeast and fruit fly we studied phosphoserines obtained in two high-throughput experiment each, by different groups of researchers [13-16]. For human we used datasets obtained in four different high-throughput experiments [17-20]. Phosphorylation is highly dynamic process, and the overlap of phosphorylation events identified in different experiments from various cell lines and tissues is relatively small. Sites observed to be phosphorylated in more than one high-throughput experiment likely are modified in a more constitutive manner, or at least represent a more reliable dataset of phosphoserines.

We analysed the evolution of modification sites and their non-modified counterparts separately among eight vertebrates (human *Homo sapiens*; chimpanzee *Pan troglodytes*; mouse *Mus musculus*; rat *Rattus norvegicus*; cow *Bos taurus*; dog *Canis lupus familiaris*; chicken *Gallus gallus*; and zebrafish *Danio rerio*), eleven fruit flies (*Drosophila melanogaster*; *D. yakuba*; *D. erecta*; *D. sechecellia*; *D. ananassae*; *D. pseudoobscura*; *D. persimilis*; *D. wilsoni*; *D. mojavensis*; *D. virilis*; *D. grimshawi*) and fifteen fungi (*Saccharomyces cerevisiae*; *S. paradoxus*; *S. mikatae*; *S. bayanus*; *Candida glabrata*; *S. castellii*; *Kluyveromyces waltii*; *K. lactis*; *Ashbya gossypii*; *Debaryomyces hansenii*; *C. albicans*; *Yarrowia lipolytica*; *Aspergillus nidulans*; *Neurospora crassa*; *Schizosaccharomyces pombe*). Orthologs of modified *H. sapiens* proteins were obtained from HomoloGene [21]; for *D. melanogaster*, from FlyBase [22]; and for *S. cerevisiae*, from FungalOrthogroups [23]. Only orthologs with the highest identity to the modified protein were selected

* Correspondence: gelfand@iitp.ru

¹Institute for Information Transmission Problems (the Kharkevich Institute) RAS, Bolshoi Karetny pereulok 19, Moscow, 127994, Russia
Full list of author information is available at the end of the article

Table 1 Datasets of phosphorylated and non-phosphorylated serines

	<i>S. cerevisiae</i>	<i>D. melanogaster</i>	<i>H. sapiens</i>
<i>Initial sets of serine residues</i>			
all phosphoserines	7381	11785	11624
phosphoserines observed more than once	1649	3137	2589
non-phosphorylated serines	103682	202574	243968
<i>Serines with at least one substitution to other types of amino acids, within ordered regions</i>			
all phosphoserines	215	180	434
phosphoserines observed more than once	21	38	43
non-phosphorylated serines	20459	13826	26350
<i>Serines with at least one substitution to other types of amino acids, within disordered regions</i>			
all phosphoserines	3666	2482	4277
phosphoserines observed more than once	857	611	906
non-phosphorylated serines	31815	42424	78120

The number of serines in each set is given. Only serines from disordered regions, with at least one substitution to other types of amino acids, were analyzed (the last three lines).

To do that that, we generated 10000 random control sets of non-phosphorylated serines. Each control set was of the same size as the corresponding phosphorylated set (generic sets and subsets of reliable phosphosites).

Structural features of phosphoserines may not be limited to disorder of surrounding protein regions, and may include other specific properties such as secondary structures, solvent availability etc. Therefore, to maximally eliminate the confounding effects, we created additional control sets containing non-modified serines located at the same protein regions as modification sites. Non-modified serines, was collected at the maximal distance of 10, 11 and 9 amino acid residues from phosphoserines, for yeast, fruit fly and human respectively. Again, the size of the control sets was the same as the size of the respective phosphoserine sets.

Differences in the substitution vectors between phosphorylated and non-phosphorylated serines from disordered regions varied among different groups of organism, but some trends were stable and significant (Figure 2). Rather unexpectedly, we did not observe any preference for substitution of phosphoserines to other aminoacids that may be phosphorylated, that is as threonine and tyrosine. At the same time, phosphorylation converts serine into a negatively charged amino acid, and, as one can see in Figure 2 in all three datasets phosphoserines are more frequently substituted to aspartate and glutamate than non-phosphorylated serines. In both cases the substitution rates of phosphoserines are much higher than in all bootstraps of control sets (P -value $\ll 10^{-4}$). In the case of the more reliable subsets of phosphoserines observed in several experiments, the substitution rate to aspartate and glutamate is even higher, and also lies outside the interval of bootstraps that in this case is wider, as the sample size is

smaller. At that, artificial substitution of serine to aspartate and glutamate, called phosphomimetic mutation, is widely used to confirm phosphorylation of serine [26,27].

There are considerable other shifts of substitution rates common to all three taxa. Particularly, phosphoserines are relatively rarely substituted to alanines and cysteines (Figure 2). However, in these cases, the control-set substitution vectors of non-phosphorylated serines located in the same regions as phosphoserines were also shifted in the same direction as phosphoserines (as compared to all non-phosphorylated serines). Hence, these shifts are likely related not to modifications, but to specific features of these regions.

The rates of substitutions to aspartate and glutamate in the additional control sets of nearest non-phosphorylated serines also are not shifted, with the exception of vertebrates where they are also shifted toward higher values (but still to a much weaker extent than in case of phosphoserines). Note that these control sets may be contaminated by phosphoserines. Indeed, phosphoserines tends to co-occur, forming clusters [28]. Therefore the sets of nearest non-phosphorylated serines likely contain phosphoserines which were not detected yet. Removing these phosphoserines would increase the significance of our observations.

The comparison with nearest non-phosphorylated serines takes into account the fact that phosphoserines tend to occur in intrinsically disordered regions. Methods used in large-scale phosphoproteomic experiments are based on selection of negatively charged peptides which results in a bias towards enrichment of phosphopeptides with acidic residues [29,30]. This fact, coupled with the fact that phosphoserines may shift positions within rapidly evolving disordered regions [31] and

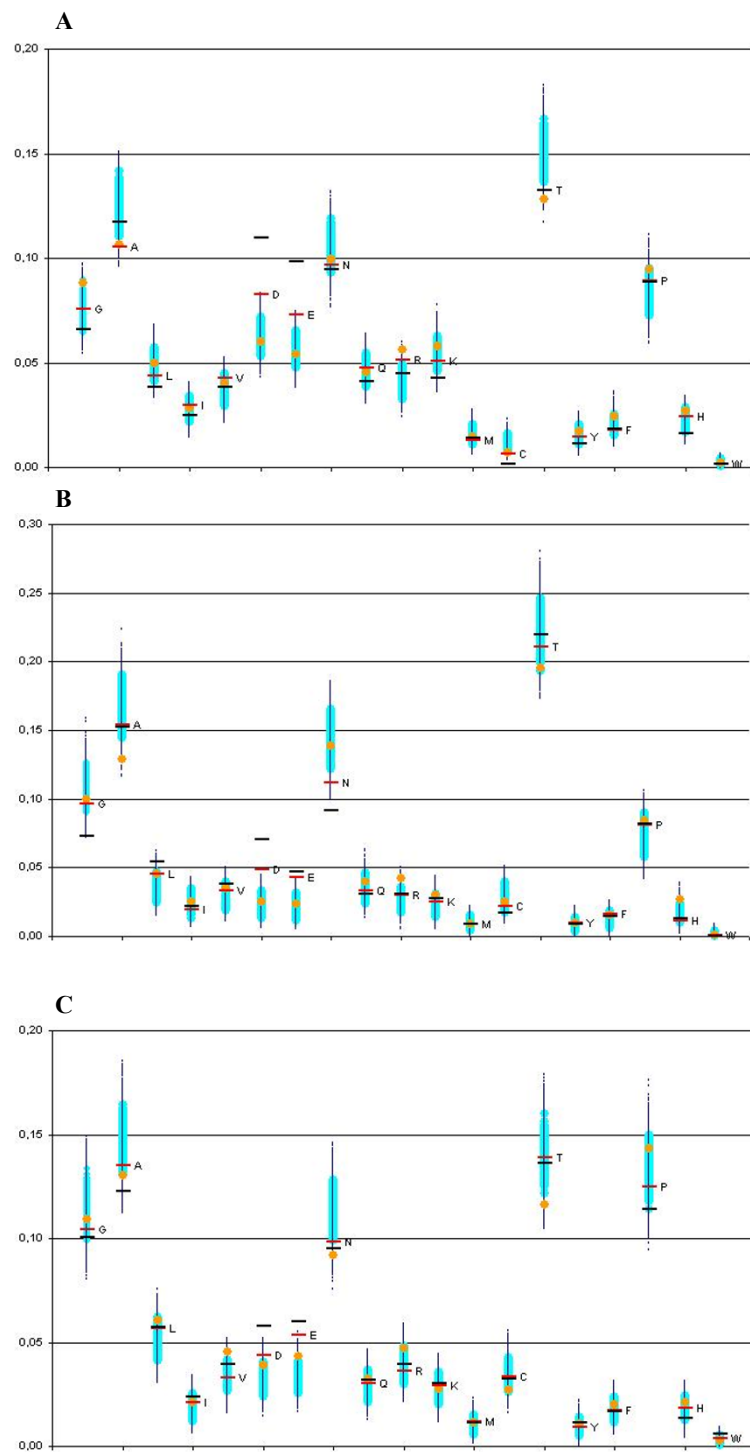


Figure 2 Substitution vectors of serines to other types of amino acids. Frequency of substitution of serines from disordered phosphoprotein regions among fungi (A), fruit flies (B) and vertebrates (C): for all phosphoserines - red bars; for phosphoserines observed in more than one experiment - black bars; 10000 control sets - clouds of large light and dark small blue dots, respectively; the additional control set of nearest non-phosphorylated serines - orange dots.

general problems of alignments of such regions could distort our analysis. But this would have the same influence on our control sets of non-modified serines from the same regions of proteins. Hence the observed differences between these controls and phosphoserines cannot be explained by such artifacts.

In addition to serine phosphorylation, we analysed the evolution of another abundant type of protein modification, lysine acetylation. Recently two large datasets of human acetylation sites became available [32,33]. We observed some differences between substitution vectors of acetylated and non-acetylated lysines, but the results obtained for these two sets of acetyllysines were discordant (data not shown). As noted in one of these papers [32], the spectrum of acetylated proteins is different between these two datasets obtained from different tissues. We observed that less than 2% of sites are common for both datasets. It seems that the available acetyllysine data are not sufficient for meaningful analysis.

It should be taken into account that our substitution vectors are probably enriched with false-positive phosphosites. This results from our over-simplified assumption that a site is modified from the first appearance of the corresponding residue in the evolutionary record. Additionally, phosphoserines from large-scale experiments may be false-positive sites. There is evidence that many phosphorylation sites could be non-functional or non-specific, as sometimes functional targets of phosphorylation are not particular sites, but entire protein regions [31,34,35]. On the other hand, the control sets could contain not yet detected phosphoserines. These false positives and false negatives should blur the differences between the substitution vectors of modified and non-modified residues. Most likely, the real level of differences is higher than the one observed here.

Reviewers' comments

Reviewer's Report 1

Reviewer 1: Arcady Mushegian - Stowers Institute, Kansas City, USA

Reviewer's comment

The idea of comparing of evolutionary substitution patterns of modified and non-modified residues in proteins is good, and the approach proposed by the authors, i.e., to reconstruct, using an ML model, the point at which the target of modification first emerged and then to see what it mutates to, is probably the only computational approach plausible at the moment.

I trust the authors that their implementation of this approach is technically sound, but, unfortunately, this is hard to ascertain from the submitted version of the manuscript, which reads as a preliminary draft devoid of

the quantitative details. This has to change - please provide at least the following:

1. The collection of phosphorylated and acetylated sites: how many sites of each type in each organism are there?

Author's Response A table with a description of the final datasets used for the construction of substitution vectors has been added to the revised version (Table 1).

Reviewer's comment

2. The phosphorylation sites at least (also acetylated sites?) are said to occur more often in the intrinsically disordered regions. Taking the non-globular regions in the proteins (which can be identified, e.g., using Wootton and Federhen's SEG program) as a proxy for "intrinsic disorder", can it be said that the actual sample of modified residues that the authors were working with is indeed more commonly occurring in such regions? And how does this sit with the ability to align the proteins in these regions?

Author's Response We predicted intrinsically disordered regions and recalculated substitution vectors separately for serine residues from disordered and ordered regions. Most of phosphoserines from the initial datasets came from protein regions predicted to be disordered (Table 1). Problems with alignments of such region are discussed in the revised version. Additional controls of non-modified serines from same regions of proteins were introduced to address this problem.

Reviewer's comment

3. The "control sets" of non-modified serines (more accurately, not-observed-to-be-modified serines): are these found in the disordered/non-globular regions to the same extent as the modified ones? If not, the controls may be biased with regard to amino acid composition and to the regions of the protein molecules (e.g., buried vs exposed) - test this directly please.

Author's Response Indeed, the amino acid composition of disordered regions and regions with a regular structure differs strongly. As described in response to comment #2, in the revised version we considered both phosphosrylated and non-phosphorylated serines from disordered and regular regions separately. Moreover, as discussed in the revised text, sets including only closest non-modified serines provide an even better control for artifacts that could be caused by specifics of regions surrounding modification sites.

Reviewer's comment

4. The trends that the authors discuss are interesting but weak - to what extent this may be explained by the small sample sizes? What was the statistical test for which the P-values are reported?

Author's Response The initial dataset of serines were large enough, but only a fraction of them were substituted to other amino acids as demonstrated by

evolutionary reconstruction. The final datasets are described in Table 1.

To measure the statistical significance, we used bootstraps of control sets of non-modified serines. For all phosphoserines and, separately, for the subset of phosphoserines observed in more than one experiment, we generated 10000 random sets of non-modified serines of appropriate size. For additional controls using neighbouring sites, we compiled sets of nearest serines of the same size as the corresponding sets of phosphorylated serines.

Reviewer's comment

5. In vertebrates, the "neighboring" serines from control set 2 seem to be faithfully following the trend towards change into D or E, with some separation from the control set 1. If this trend withstands the possible correction proposed in #2, perhaps this means that, in a "disordered" region that has several serines, any or all of them may targets of phosphorylation. Perhaps then it would be interesting to sum the substitution vectors over the region that has several serines, at least one of which is phosphorylated (i.e., how likely is it that at least one serine in this region is substituted by amino acid X?)

Author's Response The phosphoserines tends to cluster in the sequence [28]. Thus, as discussed in the revised version, the control set consisting of nearest non-phosphorylated serines could be contaminated by false-negative phosphoserines, not yet detected in experiments. On the other hand, as the trend in the control set of nearest serines is weaker, averaging of the substitution vectors would simply dilute the observation.

Reviewer's Report 2

Reviewer 2: Sandor Pongor - International Centre for Genetic Engineering and Biotechnology, Trieste, Italy

Reviewer's comment

There is mounting evidence in recent years that the study of post-translational modifications has important lessons for understanding diverse aspects of protein evolution. It has been noted among others that phosphorylated sites tend to occur in those segments of the proteins that are intrinsically disordered and/or correspond to alternative splice sites. Currently there are insufficient data on the conservation of modified sites. Kurmangalyev and associates address this problem using carefully selected datasets and well-designed statistical analyses.

The authors conclude that there are significant differences in the evolution of modified and corresponding non-modified amino acids. In particular, phosphoserines are more frequently substituted to aspartate and glutamate, compared to non-phosphorylated serines. Similarly, acetyllysines are more rarely substituted to isoleucine and

valine. These findings underline the importance of post-translational modifications when discussing the variation of residue conservations within sequence regions. The methodology is straightforward and sound and will be a useful template for future studies. The authors may want to add a few examples for situation where this approach can or can not be used.

Author's Response As discussed in the revised text, the analysis of a newly available dataset of human acetylation sites [33] did not confirm our initial observations. This is likely due to low reproducibility of currently available datasets of avetyllsines (the overlap between two datasets is extremely small). This suggests that conclusions based on such analyses should be done carefully, on data obtained from different sources and for a variety of organisms. We have encountered a similar problem with phosphothreonines and phosphotyrosines, where the datasets were simply too small for reliable conclusions.

Acknowledgements

We are grateful to Dmitry Malko, Ekaterina Ermakova and Anna Lyubetskaya who shared their programs and data, and to Stefka Tyanova and Jürgen Cox for useful discussions. This study was partially supported by the state contract 2.740.11.0101, Russian Foundation of Basic Research (09-04-92745), and program "Molecular and Cellular Biology" of the Russian Academy of Sciences.

Author details

¹Institute for Information Transmission Problems (the Kharkevich Institute) RAS, Bolshoi Karetny pereulok 19, Moscow, 127994, Russia. ²National Center for Biotechnology of the Republic of Kazakhstan, Valikhanov str., 13/1, Astana, 010000, Republic of Kazakhstan. ³Faculty of Bioengineering and Bioinformatics, Moscow State University, Vorobievyy Gory 1-73, Moscow, 119991, Russia.

Authors' contributions

YK and MG conceived the study. YK compiled the data. AG developed algorithms. YK and MG performed calculations. YK and MG analyzed the results and wrote the paper. All authors have approved the final version.

Competing interests

The authors declare that they have no competing interests.

Received: 28 September 2010 Accepted: 9 February 2011

Published: 9 February 2011

References

1. Mann M, Jensen ON: **Proteomic analysis of post-translational modifications.** *Nat Biotechnol* 2003, **21**:255-261.
2. Seo J, Lee KJ: **Post-translational Modifications and Their Biological Function: Proteomic Analysis and Systematic Approaches.** *Journal of Biochemistry and Molecular Biology* 2004, **37**:35-44.
3. Hunter T: **Signaling-2000 and beyond.** *Cell* 2000, **100**:113-127.
4. Cohen P: **The origins of protein phosphorylation.** *Nat Cell Biol* 2002, **4**: E127-E130.
5. Ptacek J, Snyder M: **Charging it up: global analysis of protein phosphorylation.** *Trends Genet* 2006, **22**:545-554.
6. Iakoucheva LM, Radivojac P, Brown CJ, O'Connor TR, Sikes JG, Obradovic Z, Dunker AK: **The importance of intrinsic disorder for protein phosphorylation.** *Nucleic Acids Res* 2004, **32**:1037-1049.
7. Gnad F, Ren S, Cox J, Olsen JV, Macek B, Oroshi M, Mann M: **PHOSIDA (phosphorylation site database): management, structural and evolutionary investigation, and prediction of phosphosites.** *Genome Biol* 2007, **8**:R250.

8. Collins MO, Yu L, Campuzano I, Grant SG, Choudhary JS: **Phosphoproteomic analysis of the mouse brain cytosol reveals a predominance of protein phosphorylation in regions of intrinsic sequence disorder.** *Mol Cell Proteomics* 2008, **7**:1331-1348.
9. Kurmangaliyev EZ, Gelfand MS: **[Alternative splicing tends to involve phosphorylation sites].** *Mol Biol (Mosk)* 2009, **43**:572-574.
10. Macek B, Gnad F, Soufi B, Kumar C, Olsen JV, Mijakovic I, Mann M: **Phosphoproteome analysis of *E. coli* reveals evolutionary conservation of bacterial Ser/Thr/Tyr phosphorylation.** *Mol Cell Proteomics* 2008, **7**:299-307.
11. Malik R, Nigg EA, Körner R: **Comparative conservation analysis of the human mitotic phosphoproteome.** *Bioinformatics* 2008, **24**:1426-1432.
12. Boekhorst J, van Breukelen B, Heck A Jr, Snel B: **Comparative phosphoproteomics reveals evolutionary and functional conservation of phosphorylation across eukaryotes.** *Genome Biol* 2008, **9**:R144.
13. Bodenmiller B, Campbell D, Gerrits B, Lam H, Jovanovic M, Picotti P, Schlapbach R, Aebersold R: **PhosphoPep - a database of protein phosphorylation sites in model organisms.** *Nat Biotechnol* 2008, **26**:1339-1340.
14. Bodenmiller B, Malmstrom J, Gerrits B, Campbell D, Lam H, Schmidt A, Rinner O, Mueller LN, Shannon PT, Pedrioli PG, Panse C, Lee HK, Schlapbach R, Aebersold R: **PhosphoPep - a phosphoproteome resource for systems biology research in *Drosophila* Kc167 cells.** *Mol Syst Biol* 2007, **3**:139.
15. Hilger M, Bonaldi T, Gnad F, Mann M: **Systems-wide analysis of a phosphataseknock-down by quantitative proteomics and phosphoproteomics.** *Mol Cell Proteomics* 2009, **8**:1908-1920.
16. Gnad F, de Godoy LM, Cox J, Neuhauser N, Ren S, Olsen JV, Mann M: **High-accuracy identification and bioinformatic analysis of in vivo protein phosphorylation sites in yeast.** *Proteomics* 2009, **9**:4642-4652.
17. Olsen JV, Blagoev B, Gnad F, Macek B, Kumar C, Mortensen P, Mann M: **Global, in vivo, and site-specific phosphorylation dynamics in signaling networks.** *Cell* 2006, **127**:635-648.
18. Daub H, Olsen JV, Bairlein M, Gnad F, Oppermann FS, Körner R, Greff Z, Kéri G, Stemmann O, Mann M: **Kinase-selective enrichment enables quantitative phosphoproteomics of the kinome across the cell cycle.** *Mol Cell* 2008, **31**:438-448.
19. Oppermann FS, Gnad F, Olsen JV, Hornberger R, Greff Z, Kéri G, Mann M, Daub H: **Large-scale proteomics analysis of the human kinome.** *Mol Cell Proteomics* 2009, **8**:1751-1764.
20. Olsen JV, Vermeulen M, Santamaria A, Kumar C, Miller ML, Jensen LJ, Gnad F, Cox J, Jensen TS, Nigg EA, Brunak S, Mann M: **Quantitative phosphoproteomics reveals widespread full phosphorylation site occupancy during mitosis.** *Sci Signal* 2010, **3**:ra3.
21. Wheeler Geer LY, Marchler-Bauer A, Geer RC, Han L, He J, He S, Liu C, Shi W, Bryant SH: **The NCBI BioSystems database.** *Nucleic Acids Res* 2010, **38**:D492-D496.
22. Tweedie S, Ashburner M, Falls K, Leyland P, McQuilton P, Marygold S, Millburn G, Osumi-Sutherland D, Schroeder A, Seal R, Zhang H, The FlyBase Consortium: **FlyBase: enhancing *Drosophila* Gene Ontology annotations.** *Nucleic Acids Res* 2009, **37**:D555-D559.
23. Wapinski I, Pfeffer A, Friedman N, Regev A: **Natural history and evolutionary principles of gene duplication in fungi.** *Nature* 2007, **449**:54-61.
24. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG: **Clustal W and Clustal X version 2.0.** *Bioinformatics* 2007, **23**:2947-2948.
25. Peng K, Radivojac P, Vucetic S, Dunker AK, Obradovic Z: **Length-dependent prediction of protein intrinsic disorder.** *BMC Bioinformatics* 2006, **7**:208.
26. Tarrant MK, Cole PA: **The chemical biology of protein phosphorylation.** *Annu Rev Biochem* 2009, **78**:797-825.
27. Song Q, Pallikkuth S, Bossuyt J, Bers DM, Robia SL: **Phosphomimetic mutations enhance phospholemmann oligomerization and modulate its interaction with the NAK-ATPase.** *J Biol Chem* 2011.
28. Schweiger R, Linal M: **Cooperativity within proximal phosphorylation sites is revealed from large-scale proteomics data.** *Biol Direct* 2010, **5**:6.
29. Ficarro SB, McClelland ML, Stukenberg PT, Burke DJ, Ross MM, Shabanowitz J, Hunt DF, White FM: **Phosphoproteome analysis by mass spectrometry and its application to *Saccharomyces cerevisiae*.** *Nat Biotechnol* 2002, **20**:301-305.
30. Mann M, Ong SE, Grønborg M, Steen H, Jensen ON, Pandey A: **Analysis of protein phosphorylation using mass spectrometry: deciphering the phosphoproteome.** *Trends Biotechnol* 2002, **20**:261-268.
31. Holt LJ, Tuch BB, Villén J, Johnson AD, Gygi SP, Morgan DO: **Global analysis of Cdk1 substrate phosphorylation sites provides insights into evolution.** *Science* 2009, **325**:1682-1686.
32. Zhao S, Xu W, Jiang W, Yu W, Lin Y, Zhang T, Yao J, Zhou L, Zeng Y, Li H, Li Y, Shi J, An W, Hancock SM, He F, Qin L, Chin J, Yang P, Chen X, Lei Q, Xiong Y, Guan K: **Regulation of Cellular Metabolism by Protein Lysine Acetylation.** *Science* 2010, **327**:1000-1004.
33. Choudhary C, Kumar C, Gnad F, Nielsen ML, Rehman M, Walther TC, Olsen JV, Mann M: **Lysine acetylation targets protein complexes and co-regulates major cellular functions.** *Science* 2009, **325**:834-840.
34. Landry CR, Levy ED, Michnick SW: **Weak functional constraints on phosphoproteomes.** *Trends Genet* 2009, **25**:193-197.
35. Tan CS, Jørgensen C, Linding R: **Roles of "junk phosphorylation" in modulating biomolecular association of phosphorylated proteins?** *Cell Cycle* 2010, **9**:1276-1280.

doi:10.1186/1745-6150-6-8

Cite this article as: Kurmangaliyev et al.: Evolutionary patterns of phosphorylated serines. *Biology Direct* 2011 6:8.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

