

CPLM: a database of protein lysine modifications

Zexian Liu¹, Yongbo Wang¹, Tianshun Gao¹, Zhicheng Pan¹, Han Cheng¹, Qing Yang¹, Zhongyi Cheng², Anyuan Guo¹, Jian Ren³ and Yu Xue^{1,*}

¹Department of Biomedical Engineering, College of Life Science and Technology, Huazhong University of Science and Technology, Wuhan, Hubei 430074, China, ²Advanced Institute of Translational Medicine, Tongji University, Shanghai 200092, China and ³State Key Laboratory of Biocontrol, School of Life Sciences, Sun Yat-sen University, Guangzhou, Guangdong 510275, China

Received September 13, 2013; Revised and Accepted October 17, 2013

ABSTRACT

We reported an integrated database of Compendium of Protein Lysine Modifications (CPLM; <http://cplm.biocuckoo.org>) for protein lysine modifications (PLMs), which occur at active ϵ -amino groups of specific lysine residues in proteins and are critical for orchestrating various biological processes. The CPLM database was updated from our previously developed database of Compendium of Protein Lysine Acetylation (CPLA), which contained 7151 lysine acetylation sites in 3311 proteins. Here, we manually collected experimentally identified substrates and sites for 12 types of PLMs, including acetylation, ubiquitination, sumoylation, methylation, butyrylation, crotonylation, glycation, malonylation, phosphoglyceration, propionylation, succinylation and pupylation. In total, the CPLM database contained 203972 modification events on 189919 modified lysines in 45748 proteins for 122 species. With the dataset, we totally identified 76 types of co-occurrences of various PLMs on the same lysine residues, and the most abundant PLM crosstalk is between acetylation and ubiquitination. Up to 53.5% of acetylation and 33.1% of ubiquitination events co-occur at 10746 lysine sites. Thus, the various PLM crosstalks suggested that a considerable proportion of lysines were competitively and dynamically regulated in a complicated manner. Taken together, the CPLM database can serve as a useful resource for further research of PLMs.

INTRODUCTION

In 1964, Allfrey *et al.* (1) first observed gene expression regulation mediated by covalently introducing acetyl and methyl groups on lysine residues in histones. Numerous following studies in epigenetics proposed the

combinational post-translational modifications (PTMs) of histones as ‘histone codes’, of which PTMs occurring on lysine residues occupy an important proportion (2). Later studies discovered lysine as a hot spot for PTMs, while a number of protein lysine modifications (PLMs) can occur in both histone and non-histone proteins (3–11). For example, beyond constituting the ‘histone code’, lysine acetylation plays a critical role in various biological processes such as metabolism (12,13) and autophagy (14,15), while methylation in non-histone proteins can regulate protein stability and activity (16). In 2004, the Nobel Prize in Chemistry was awarded to Aaron Ciechanover, Avram Herschko and Irwin Rose for their discovery of ubiquitin conjugation on lysine as a mechanism that targets proteins for degradation (17). Also, ubiquitin-like proteins such as small ubiquitin-related modifier and prokaryotic ubiquitin-like protein were found to modify protein lysine residues through a conserved conjugation cascade (18,19). In addition, protein lysines can be modified to 3-phosphoglyceryl-lysine by the primary glycolytic intermediate 1,3-bisphosphoglycerate (1,3-BPG) (10), whereas lysine glycation is involved in glycolytic processes (11).

Recently, rapid progresses in proteomic technologies greatly advanced the identification of well-characterized PLMs (20–23) and the discovery of new PLMs (4,6–8,10). For example, with a monoclonal antibody for diglycine (diGly)-containing isopeptides, Kim *et al.* (21) identified and quantified nearly 20 000 ubiquitination sites. Also, Udeshi *et al.* (22) refined a preparation procedure and used anti-diGly antibodies to quantify ~20 000 ubiquitination sites. In 2012, Lundby *et al.* (23) quantified ~15 000 acetylation sites from 16 rat tissues and systematically analyzed the tissue-specific lysine acetylation profiles. In particular, with the state-of-the-art proteomic techniques, Dr. Yingming Zhao’s group has identified a number of new PLMs such as butyrylation (4), propionylation (4), malonylation (6), crotonylation (7) and succinylation (8). Because the numbers of PLMs and modified lysine residues have been greatly expanded,

*To whom correspondence should be addressed. Tel: +86 27 87793903; Fax: +86 27 87793172; Email: xueyu@hust.edu.cn

an integrated resource for the community is urgently needed. Although several public databases such as UniProt (24), HPRD (25), SysPTM (26) and dbPTM (27) contained information for PLMs, only a limited proportion of the identified substrates and sites were covered, and the newly discovered PLMs were not considered.

Previously, we developed the Compendium of Protein Lysine Acetylation (CPLA) database to maintain the identified lysine acetylation information (28). In this work, we greatly improved the CPLA database by extending the types of PLMs and developed the database of Compendium of Protein Lysine Modifications (CPLM). From scientific literature, the experimentally identified substrates and sites for 12 types of PLMs were manually collected. Besides acetylation, well-studied PLMs such as ubiquitination, sumoylation, methylation and glycation and newly discovered PLMs including butyrylation, crotonylation, malonylation, phosphoglycerylation, propionylation, succinylation and pupylation were integrated into the database. Currently, CPLM database contained 203 972 modification events on 189 919 modified lysine residues in 45 748 proteins from 122 species, and the detailed annotations were also provided. The database can be searched or browsed in a convenient manner. Based on the comprehensive dataset, we systematically analyzed the concurrences of different PLMs at the same lysine residues. Although the number of identified substrates and sites for different types of PLMs varies from ten thousands to tens, each PLM can crosstalk with at least one other PLM and the co-occurrences of different PLMs at the same site were particularly abundant. From 76 types of identified PLM co-occurrences, we observed that the crosstalks among acetylation, ubiquitination and succinylation are mostly abundant. The intensive crosstalks among PLMs suggested that at least a considerable number of lysines were competitively and dynamically regulated by different PLMs. Taken together, the CPLM database provided an integrative platform for the community to access the current processes on PLMs and generated a useful resource for further experimental or computational considerations. The CPLM database was implemented in PHP + MySQL + JavaScript.

CONSTRUCTION AND CONTENT

As previously described (28), we searched PubMed with keywords including ‘acetylation’, ‘ubiquitination’, ‘sumoylation’, ‘methylation’, ‘glycation’, ‘butyrylation’, ‘crotonylation’, ‘malonylation’, ‘phosphoglycerylation’, ‘propionylation’, ‘succinylation’ and ‘pupylation’ and manually curated literature to collect the experimentally identified PLM substrates and sites. To avoid missing data, additional keywords such as ‘acetylated’, ‘acetyl’, ‘ubiquitinated’ and other related nomenclatures were employed for searching more data in PubMed. All modified lysine residues were mapped to the benchmark sequences retrieved from the UniProt database (Release 2013_08) (24). To provide more information for the PTMs substrates, the annotations from UniProt (24)

were integrated into the database. The primary references for PLM substrates and sites were also provided to ensure the quality of the database.

In total, 203 972 modification events were found to occur on 189 919 lysine residues in 45 748 substrates for 12 types of PLMs (Supplementary Table S1). Obviously, acetylation and ubiquitination have the most substrates; the former contains 58 563 sites in 20 088 proteins and the latter contains 139 950 sites in 32 429 proteins (Supplementary Table S1). The third PLM with most substrates is succinylation (8), which was discovered as a novel PLM in 2011 and identified with 2523 sites in 897 substrates (Supplementary Table S1). The rapid progress in the identification of succinylation is attributed to the advancement of proteomic techniques (29). However, for other new PLMs such as butyrylation, crotonylation, malonylation, phosphoglycerylation and propionylation, there were only a small number of identified substrates that mainly focused on histones (Supplementary Table S1). Although various PLMs were experimentally detected in 122 species, the number of identified substrates is usually limited for most organisms. With the ggplot2 program (30) in the R package (31), the distribution of PLM substrates and sites from 12 major species with >200 substrates were visualized (Figure 1A and B). Clearly, animals, especially mammals, were identified with most substrates (Figure 1A) and sites (Figure 1B). It is worthy to note that several types of PLMs are only exclusively identified in distinct species. For example, ubiquitination and sumoylation are only available in eukaryotes, while pupylation was only discovered in actinomycetes.

USAGE

The CPLM database was developed in a user-friendly manner, while browse and search options were provided for accessing the information. Because the proteins and sites could be classified according to the PLM types and species, two browse options including ‘Browse by types’ and ‘Browse by species’ were developed in the database (Figure 2). For convenience, only 12 major species were listed for browsing, while all the other organisms were denoted as ‘Others’. Here, we use lysine acetylation substrates from *Homo sapiens* as an example to present the usage of the browse options in CPLM. In the option of ‘Browse by types’, 12 simplified molecular structures of ligands conjugated to lysine residues during modification were employed to represent the 12 types of PLMs (Figure 2A). By clicking on the ‘Acetylation’ button, a brief introduction of protein lysine acetylation and the protein number distribution of acetylated proteins in 12 major organisms and other species were showed (Figure 2A). Then the acetylation substrates in *H. sapiens* could be listed through clicking on the ‘*Homo sapiens*’ link (Figure 2B). In the option of ‘Browse by species’, the 12 major organisms were organized as animals, bacteria, fungi and plants. Users could click on the ‘*H. sapiens*’ button to view the protein number distribution of different PLM substrates in *H. sapiens* (Figure 2C), and then click on the link of ‘Acetylation’

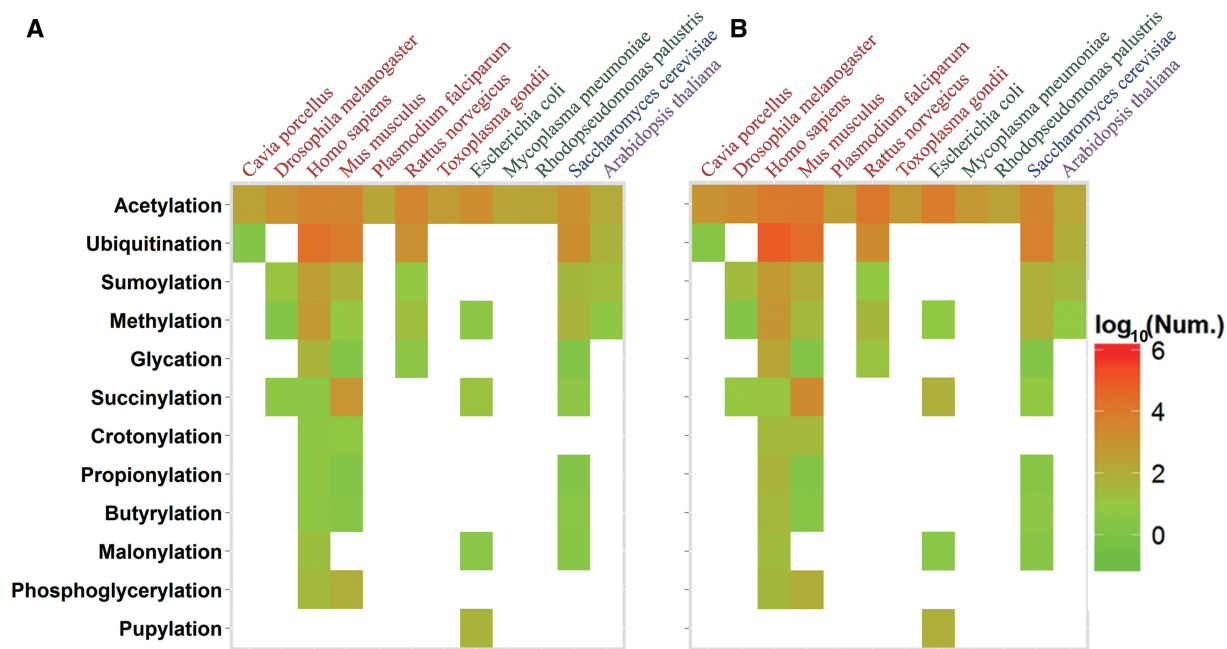


Figure 1. The heatmaps for the protein number distribution of different PLM types and species. The species names in red, green, blue and purple are from animals, bacteria, fungi and plants, respectively. (A) The heatmap for the number of substrates; (B) the heatmap for the number of modified lysine residues.

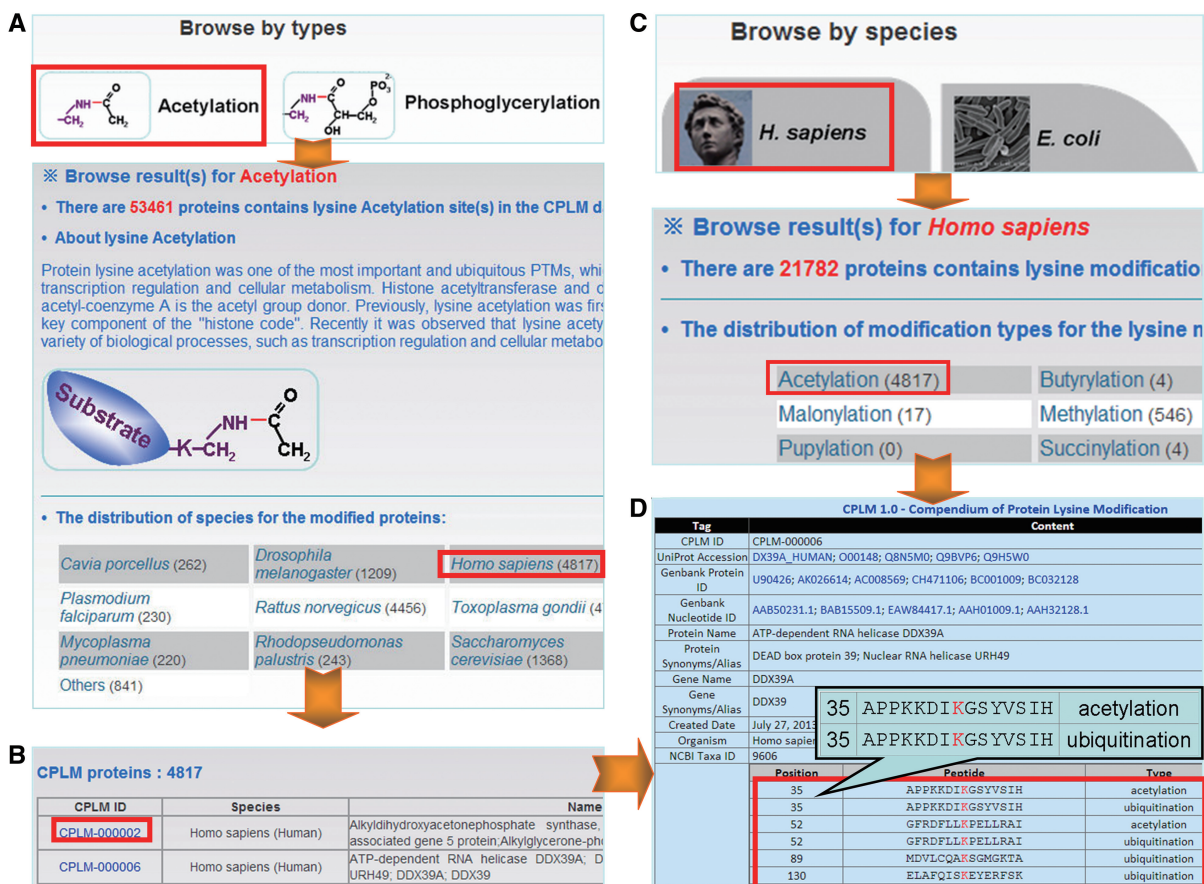


Figure 2. The browse options of CPLM. Two browse approaches including by PLM types and by species were provided to browse the database. (A) By PLM types; (B) the protein list for specified PLM and selected organism; (C) by species; (D) the detailed information of human dead box protein 39.

to view the list of acetylated substrates in *H. sapiens* (Figure 2B). The detailed information for any specified protein could be accessed through the links in the list (Figure 2D).

For convenient usage, three search options were implemented for querying the database with one or multiple keywords. For example, if users search the keyword 'TP53' in the 'Gene Name' area, the results will be shown in a tabular format with CPLM ID, organism and protein/gene names/aliases (Figure 3A). Furthermore, two options including 'Advance Search' and 'BLAST Search' were developed to query the proteins with higher accuracy. In the 'Advance Search' option, users can submit up to three search terms, which could be specified in different areas and combined with three operators of 'and', 'or' and 'exclude' to perform a complex query (Figure 3B). The 'BLAST search' option was designed to find similar proteins with a protein sequence in the FASTA format. Through the application of NCBI BLAST packages (32), users could submit a protein sequence in the FASTA format to search identical or homologous proteins (Figure 3C).

DISCUSSION

As an important molecular mechanism, PTMs greatly expand the proteome complexity and play a critical role

in the regulation of various biological processes (20,33). With the active ε-amino groups, lysine residues were modified by various PLMs, which constitute an important proportion among the large number of PTM types (3). Through modifying the substrates, PLMs regulate various biological processes, while aberrances of lysine modifications were associated with diseases and cancers (18,34–36). Recent development of proteomic techniques greatly advances the identification of PLM substrates and the discovery of new types of PLMs (3,20). However, in contrast to other PTMs such as phosphorylation (37,38), the computational resources for PLMs are still limited.

In this work, we updated the acetylation-associated database of CPLA into CPLM for more types of PLMs. Because 203 972 modification events for 12 types of PLMs were identified on 189 919 lysine residues, it was expected that there were a large number of co-occurrences among different PTMs. Indeed, Weinert *et al.* (39) discovered that the crosstalks between acetylation and succinylation are extensive in both prokaryotes and eukaryotes. Also, previous studies identified that the competition between acetylation and ubiquitination can serve as a mechanism to control protein stability (40) and activity (41). From the data set, we totally identified 76 types of PLM co-occurrences at same lysine residues, including 40 types of pairwise crosstalks (Figure 4A) and 36 types of multiple

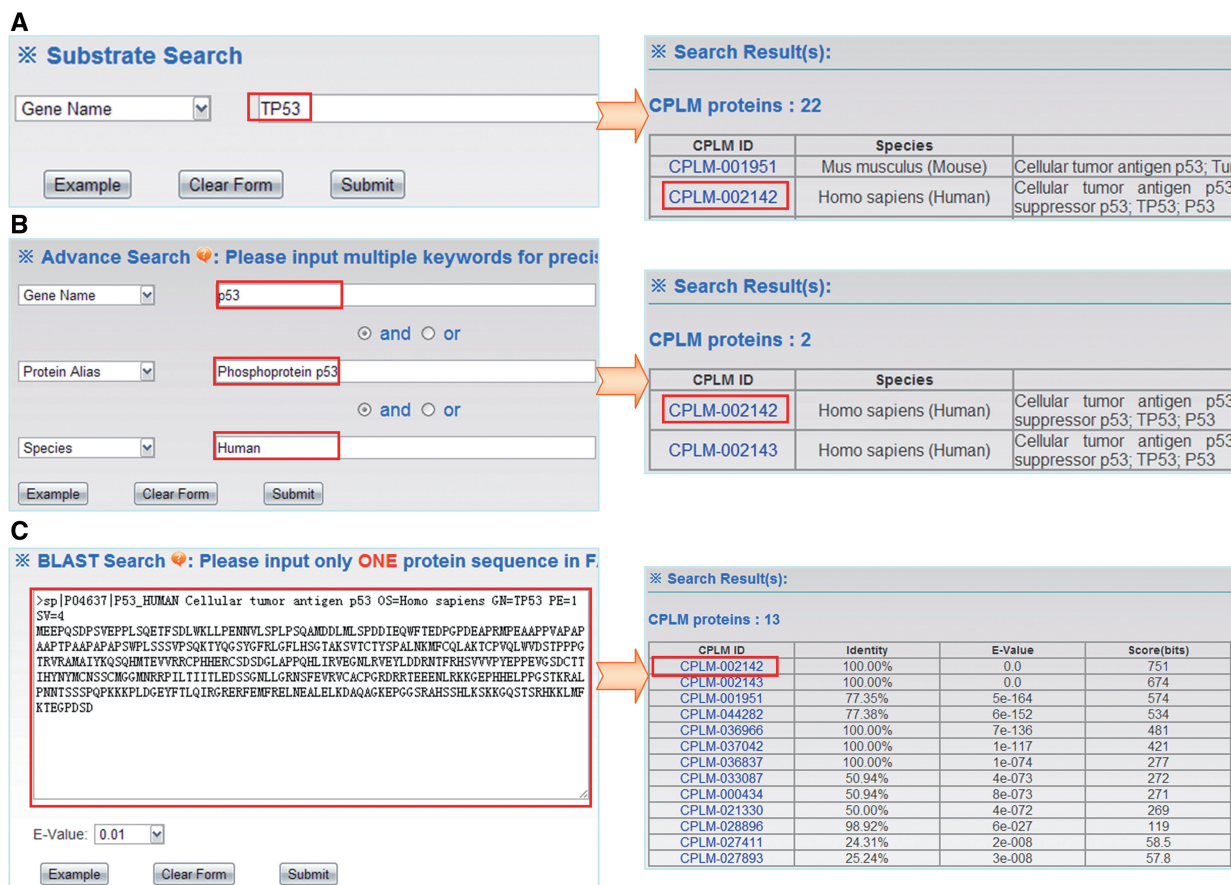


Figure 3. The search options. (A) The database could be queried with simple keywords input; (B) the 'Advance Search' allows users to submit combination of up to three terms for searching; (C) the database could be queried with a protein sequence to find identical or homologous proteins.

Science & Technology Cooperation Program of China [OS2013ZR0003].

Conflict of interest statement. None declared.

REFERENCES

- Allfrey, V.G., Faulkner, R. and Mirsky, A.E. (1964) Acetylation and methylation of histones and their possible role in the regulation of RNA synthesis. *Proc. Natl Acad. Sci. USA*, **51**, 786–794.
- Jenuwein, T. and Allis, C.D. (2001) Translating the histone code. *Science*, **293**, 1074–1080.
- Olsen, C.A. (2012) Expansion of the lysine acylation landscape. *Angew. Chem. Int. Ed. Engl.*, **51**, 3755–3756.
- Chen, Y., Sprung, R., Tang, Y., Ball, H., Sangras, B., Kim, S.C., Falck, J.R., Peng, J., Gu, W. and Zhao, Y. (2007) Lysine propionylation and butyrylation are novel post-translational modifications in histones. *Mol. Cell Proteomics*, **6**, 812–819.
- Cheng, Z., Tang, Y., Chen, Y., Kim, S., Liu, H., Li, S.S., Gu, W. and Zhao, Y. (2009) Molecular characterization of propionyllysines in non-histone proteins. *Mol. Cell Proteomics*, **8**, 45–52.
- Peng, C., Lu, Z., Xie, Z., Cheng, Z., Chen, Y., Tan, M., Luo, H., Zhang, Y., He, W., Yang, K. *et al.* (2011) The first identification of lysine malonylation substrates and its regulatory enzyme. *Mol. Cell Proteomics*, **10**, M111.012658.
- Tan, M., Luo, H., Lee, S., Jin, F., Yang, J.S., Montellier, E., Buchou, T., Cheng, Z., Rousseaux, S., Rajagopal, N. *et al.* (2011) Identification of 67 histone marks and histone lysine crotonylation as a new type of histone modification. *Cell*, **146**, 1016–1028.
- Zhang, Z., Tan, M., Xie, Z., Dai, L., Chen, Y. and Zhao, Y. (2011) Identification of lysine succinylation as a new post-translational modification. *Nat. Chem. Biol.*, **7**, 58–63.
- Xie, Z., Dai, J., Dai, L., Tan, M., Cheng, Z., Wu, Y., Boeke, J.D. and Zhao, Y. (2012) Lysine succinylation and lysine malonylation in histones. *Mol. Cell Proteomics*, **11**, 100–107.
- Moellering, R.E. and Cravatt, B.F. (2013) Functional lysine modification by an intrinsically reactive primary glycolytic metabolite. *Science*, **341**, 549–553.
- Ansari, N.A., Moinuddin, and Ali, R. (2011) Glycated lysine residues: a marker for non-enzymatic protein glycation in age-related diseases. *Dis. Markers*, **30**, 317–324.
- Zhao, S., Xu, W., Jiang, W., Yu, W., Lin, Y., Zhang, T., Yao, J., Zhou, L., Zeng, Y., Li, H. *et al.* (2010) Regulation of cellular metabolism by protein lysine acetylation. *Science*, **327**, 1000–1004.
- Wang, Q., Zhang, Y., Yang, C., Xiong, H., Lin, Y., Yao, J., Li, H., Xie, L., Zhao, W., Yao, Y. *et al.* (2010) Acetylation of metabolic enzymes coordinates carbon source utilization and metabolic flux. *Science*, **327**, 1004–1007.
- Yi, C., Ma, M., Ran, L., Zheng, J., Tong, J., Zhu, J., Ma, C., Sun, Y., Zhang, S., Feng, W. *et al.* (2012) Function and molecular mechanism of acetylation in autophagy regulation. *Science*, **336**, 474–477.
- Lin, S.Y., Li, T.Y., Liu, Q., Zhang, C., Li, X., Chen, Y., Zhang, S.M., Lian, G., Liu, Q., Ruan, K. *et al.* (2012) GSK3-TIP60-ULK1 signaling pathway links growth factor deprivation to autophagy. *Science*, **336**, 477–481.
- Huang, J. and Berger, S.L. (2008) The emerging field of dynamic lysine methylation of non-histone proteins. *Curr. Opin. Genet. Dev.*, **18**, 152–158.
- Hershko, A. (2005) The ubiquitin system for protein degradation and some of its roles in the control of the cell-division cycle (Nobel lecture). *Angew. Chem. Int. Ed. Engl.*, **44**, 5932–5943.
- Johnson, E.S. (2004) Protein modification by SUMO. *Annu. Rev. Biochem.*, **73**, 355–382.
- Pearce, M.J., Mintseris, J., Ferreyra, J., Gygi, S.P. and Darwin, K.H. (2008) Ubiquitin-like protein involved in the proteasome pathway of *Mycobacterium tuberculosis*. *Science*, **322**, 1104–1107.
- Choudhary, C. and Mann, M. (2010) Decoding signalling networks by mass spectrometry-based proteomics. *Nat. Rev. Mol. Cell Biol.*, **11**, 427–439.
- Kim, W., Bennett, E.J., Huttlin, E.L., Guo, A., Li, J., Possemato, A., Sowa, M.E., Rad, R., Rush, J., Comb, M.J. *et al.* (2011) Systematic and quantitative assessment of the ubiquitin-modified proteome. *Mol. Cell*, **44**, 325–340.
- Udeshi, N.D., Svinkina, T., Mertins, P., Kuhn, E., Mani, D.R., Qiao, J.W. and Carr, S.A. (2013) Refined preparation and use of anti-diglycine remnant (K-epsilon-GG) antibody enables routine quantification of 1000s of ubiquitination sites in single proteomics experiments. *Mol. Cell Proteomics*, **12**, 825–831.
- Lundby, A., Lage, K., Weinert, B.T., Bekker-Jensen, D.B., Secher, A., Skovgaard, T., Kelstrup, C.D., Dmytryiev, A., Choudhary, C., Lundby, C. *et al.* (2012) Proteomic analysis of lysine acetylation sites in rat tissues reveals organ specificity and subcellular patterns. *Cell Rep.*, **2**, 419–431.
- The UniProt Consortium. (2013) Update on activities at the universal protein resource (UniProt) in 2013. *Nucleic Acids Res.*, **41**, D43–D47.
- Keshava-Prasad, T.S., Goel, R., Kandasamy, K., Keerthikumar, S., Kumar, S., Mathivanan, S., Telikicherla, D., Raju, R., Shafreen, B., Venugopal, A. *et al.* (2009) Human protein reference database—2009 update. *Nucleic Acids Res.*, **37**, D767–D772.
- Li, H., Xing, X., Ding, G., Li, Q., Wang, C., Xie, L., Zeng, R. and Li, Y. (2009) SysPTM: a systematic resource for proteomic research on post-translational modifications. *Mol. Cell. Proteomics*, **8**, 1839–1849.
- Lu, C.T., Huang, K.Y., Su, M.G., Lee, T.Y., Bretana, N.A., Chang, W.C., Chen, Y.J., Chen, Y.J. and Huang, H.D. (2013) DbPTM 3.0: an informative resource for investigating substrate site specificity and functional association of protein post-translational modifications. *Nucleic Acids Res.*, **41**, D295–D305.
- Liu, Z., Cao, J., Gao, X., Zhou, Y., Wen, L., Yang, X., Yao, X., Ren, J. and Xue, Y. (2011) CPLA 1.0: an integrated database of protein lysine acetylation. *Nucleic Acids Res.*, **39**, D1029–D1034.
- Park, J., Chen, Y., Tishkoff, D.X., Peng, C., Tan, M., Dai, L., Xie, Z., Zhang, Y., Zwaans, B.M., Skinner, M.E. *et al.* (2013) SIRT5-mediated lysine desuccinylation impacts diverse metabolic pathways. *Mol. Cell*, **50**, 919–930.
- Wickham, H. (2009) *ggplot2: Elegant Graphics For Data Analysis*. Springer, New York.
- Team, R.C. (2013) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Boratyn, G.M., Camacho, C., Cooper, P.S., Coulouris, G., Fong, A., Ma, N., Madden, T.L., Matten, W.T., McGinnis, S.D., Merezuk, Y. *et al.* (2013) BLAST: a more efficient report with usability improvements. *Nucleic Acids Res.*, **41**, W29–W33.
- Xue, Y., Liu, Z., Cao, J. and Ren, J. (2011) *Bioinformatics - Experimental Biology Systems*. InTech, Rijeka, Croatia.
- Yang, X.J. (2004) The diverse superfamily of lysine acetyltransferases and their roles in leukemia and other diseases. *Nucleic Acids Res.*, **32**, 959–976.
- Yang, X.J. and Seto, E. (2007) HATs and HDACs: from structure, function and regulation to novel strategies for therapy and prevention. *Oncogene*, **26**, 5310–5318.
- Bedford, L., Lowe, J., Dick, L.R., Mayer, R.J. and Brownell, J.E. (2011) Ubiquitin-like protein conjugation and the ubiquitin-proteasome system as drug targets. *Nat. Rev. Drug Discov.*, **10**, 29–46.
- Liu, Z., Wang, Y. and Xue, Y. (2013) Phosphoproteomics-based network medicine. *FEBS J.*, **280**, 5696–5704.
- Xue, Y., Gao, X., Cao, J., Liu, Z., Jin, C., Wen, L., Yao, X. and Ren, J. (2010) A summary of computational resources for protein phosphorylation. *Curr. Protein Pept. Sci.*, **11**, 485–496.
- Weinert, B.T., Scholz, C., Wagner, S.A., Iesmantavicius, V., Su, D., Daniel, J.A. and Choudhary, C. (2013) Lysine succinylation is a frequently occurring modification in prokaryotes and eukaryotes and extensively overlaps with acetylation. *Cell Rep.*, **4**, 842–851.
- Gronroos, E., Hellman, U., Heldin, C.H. and Ericsson, J. (2002) Control of Smad7 stability by competition between acetylation and ubiquitination. *Mol. Cell*, **10**, 483–493.
- Li, H., Wittwer, T., Weber, A., Schneider, H., Moreno, R., Maine, G.N., Kracht, M., Schmitz, M.L. and Burstein, E. (2012) Regulation of NF-kappaB activity by competition between RelA acetylation and ubiquitination. *Oncogene*, **31**, 611–623.