

元智大學

資訊工程學系

碩士論文

硫端棕櫚酰化位置於肝組織具有受質專一性之辨識調查

**Investigation and identification of S-palmitoylation sites with
substrate specificity in the liver tissues.**

研究生： Ho Thi Trang

指導教授： Dr. Tzu-Ya Weng

共同指導教授: Dr. Tzong-Yi Lee

中華民國 一百零五年 一月

硫端棕櫚酰化位置於肝組織具有受質專一性之辨識調查
Investigation and identification of S-palmitoylation sites with
substrate specificity in the liver tissues.

研究生:	胡氏妝	Student:	HO THI TRANG
指導教授:	翁資雅博士	Advisor:	Dr. TZU-YA WENG
共同指導教授:	李宗夷博士	Joint advisor:	Dr. TZONG-YI LEE

元 智 大 學

資 訊 工 程 學 系

碩 士 論 文

A Thesis
Submitted to Department of Computer Science & Engineering
Yuan Ze University
in Partial Fulfillment of the Requirements
for the Degree of Master
in
Computer Science and Engineering
January 2016
Chungli, Taiwan, Republic of China.

中 華 民 國 一 百 零 五 年 一 月

硫端棕櫚酰化位置於肝組織具有受質專一性之辨識調查

學生： 胡氏啟

指導教授： 翁育雅 博士

共同指導教授： 李宗炎 博士

元智大學資訊工程學系

摘要

硫端棕櫚酰化是經由 酶的催化下飽和的 16 個碳的棕櫚酸鹽通過硫酯鍵共價修飾到蛋白質半胱氨酸的巯基上硫鍵共價連接，是轉譯後修飾的形式之一。它是動態、可逆的，在蛋白質-蛋白質交互作用與蛋白質運輸過程中息息相關，也涉及到生物過程的多樣性...等等，是非常重要的蛋白質轉譯後修飾之一。而現今對於硫端棕櫚酰化位置受質專一性的辨識還沒有找到一個有效的方法去辨識。因此運用資訊技術找尋有效的計算方法成爲了十分迫切的需求，希望能以資訊的方法快速辨別硫端棕櫚酰化位置。本研究於 UniProtKB 資料庫以及兩篇文獻(Forrester MT et al; Yang W et al)中蒐集了肝組織中硫端棕櫚酰化位置的實驗驗證資料，總共整理出 710 個位置資料。並提出了一種遞迴統計的方法來辨識肝組織中硫端棕櫚酰化位置受質之專一性。並以五倍交叉驗證之形式使用 SVM 分類器，進行模擬測試及評估效能。結果表示所開發的方法的確有效的提升敏感性，特異性和準確性。另外使用獨立測試資料對硫端棕櫚酰化位置測試之結果進行評估。最後，本研究將此新穎的方法研究結果，呈現於免費的網站，並可幫助生物學家分析硫端棕櫚酰化位置。(網址:

<http://csb.cse.yzu.edu.tw/MDDPalm/>)

關鍵字：後轉譯修飾作用, 最大依賴分解度.

Investigation and identification of S-palmitoylation sites with substrate specificity in the liver tissues.

Student : Ho Thi Trang

Advisor : Dr. Tzu-Ya Weng
Joint advisor: Dr. Tzong-Yi Lee

Department of Computer Science and Engineering
Yuan Ze University

Abstract

S-palmitoylation, the covalent attachment of 16-carbon palmitic acids to a cysteine residue via a thioester linkage, is an important reversible lipid modification (PTM) that regulates protein trafficking, protein-protein interaction. It also related to diversity of physiological and biological process, etc... However, the substrate specificity of cysteine S-palmitoylation remains unknown. Thus, finding the effective computational method to predict S-palmitoylation sites is urgent demand in bioinformatics. Based on total of 710 experimentally verified S-palmitoylation sites in the liver tissues were collected from UniProtKB, Forrester MT et al and Yang W et al papers. This study presents a recursively statistical method to identify conserved substrate motifs for S-palmitoylation in the liver tissues. Statistical significance Support vector machine (SVM) was applied to construct predictive model learned from verified substrate motifs. The evaluation of five-fold cross-validation indicated that the model trained with identified motifs were effective in identification of S-palmitoylation sites with an enhanced sensitivity, specificity, and accuracy. It also provided a promising performance in an independent testing set. The correct identification of previous report S-palmitoylation sites of mouse protein demonstrated the effectiveness of the proposed method and it indicated that the proposed method could be a practicable ways of conducting primary analyses of proteins S-palmitoylation in the liver tissues. Finally, the constructed models have been implemented as a web-based system freely available at <http://csb.cse.yzu.edu.tw/MDDPalm/> for identifying uncharacterized S-palmitoylation sites on the protein sequences.

Keywords: S-palmitoylation, Support vector machine, maximal dependence decomposition.

Acknowledgement

Firstly, I would like to express my thanks to Yuan Ze University for the full scholarship which has illuminated my long-lasting dream of obtaining a Master Degree in Computer Science and Engineering.

Second, I am thankful to my thesis advisors, Professor Tzong-Yi Lee, Professor Tzu-Ya Weng who suggested an interesting research topic to me and helped to keep it on the right track.

Finally, my sincere thanks to my family, my best friend, my labmates, all the YZU teachers and friends that I have ever met during the two years of study in Taiwan. This thesis could be finished on time is largely due to their support.



Contents

摘要	iii
Abstract	iv
Acknowledgement	v
Contents	vi
List of Tables	viii
List of Figures	ix
Chapter 1. Introduction	1
1.1 Background	1
1.1.1 Posttranslational modification	1
1.1.2 Protein S-palmitoylation	2
1.2 Related work	3
1.3 Motivation and Goal	4
Chapter 2. Material And Methods	6
2.1 Flow Chart	6
2.2 UniprotKB/Swiss-Prot Database	6
2.3 Data Collection and Processing	7
2.4 Feature Investigation	9
2.4.1 Amino Acid (AA)	9
2.4.2 Amino Acid Composition (AAC)	10
2.4.3 Amino Acid Pair Composition (AAPC)	10
2.4.4 Blocks Substitution Matrix (Blosum62)	11
2.4.5 Accessible Surface Area (ASA)	12
2.4.6 positional weight matrix (PWM)	13
2.4.7 Position Specific Scoring Matrix (PSSM)	14
2.5 Support Vector Machine (SVM)	15
2.6 Data clustering by maximal dependence decomposition (MDDLogo)	15
2.7 Model Training and cross validation evaluation	17
2.8 Performance measurement	18
Chapter 3. Result	20
3.1 Amino acid composition analysis	20
3.2 Five-fold cross-validation performance of training features	23
3.3 MDD-clustered substrate motifs for S-palmitoylation sites and the cross-	

validation performances	24
3.4 Independent testing comparison with existing prediction tools	27
3.5 Comparison with existing prediction tools	28
3.6 Implementation of Web Server for the identification of S-palmitoylation sites ..	29
Chapter 4. Conclusion	32
4.1 Research Contributions	32
4.2 Limitation and future work	33



List of Tables

Table 1. Statistical table of palmitoylation prediction tools.	3
Table 2. Data statistics of S-palmitoylated protein on tissues of training data set and Uniprot.	5
Table 3. Data statistics of the liver tissues and the other tissues in training data set.	8
Table 4. Data statistics of training set and independent testing set.	9
Table 5. Five-fold cross validation results on single SVM model trained with various features.	24
Table 6. The 5 MDDLogo-clustered subgroups and their performances of five-fold cross-validation from 710 S-palmitoylation sites in the liver tissues of training data set.	26
Table 7. Performance of CSS-Palm 4.0, NBA-Palm, CSKAAP-Palm, WAP-Palm, PalmPred, SeqPalm and MDDPalm on the independent dataset of 566 proteins.	28



List of Figures

Figure 1. Types of Post-translation modifications.....	1
Figure 2. Structure of Protein S-palmitoylation.....	2
Figure 3. System flow of this study.....	6
Figure 4. Uniprot release 2015_09.....	7
Figure 5. Illustration of Amino Acid Composition feature.....	10
Figure 6. Illustration of Amino Acid Pair Composition feature.....	11
Figure 7. Substitution matrix for BLOSUM62.....	11
Figure 8. Illustration of the solvent accessible surface in comparison to the van der Waals surface.....	12
Figure 9. Illustration of positional weight matrix.....	13
Figure 10. Illustration of the PSSM features generation from Original PSSM profiles.....	14
Figure 11. The principle of support vector machine.....	15
Figure 12. The maximal dependence of position is detected by chi-square testing.....	16
Figure 13. Illustration of the five-fold cross-validation.....	18
Figure 14. diagram of compositional biases of amino acids around S-palmitoylation sites compared to the non-S-palmitoylation sites in the liver tissues.....	21
Figure 15. Sequence logo of S-palmitoylation sites (generated by WebLogo server).....	21
Figure 16. Sequence logo of non-S-palmitoylation sites (generated by WebLogo server).....	22
Figure 17. TwoSampleLogo presents the compositional biases of amino acids around S-palmitoylation sites compared to the non-S-palmitoylation sites in the liver tissues.....	22
Figure 18. Performance comparison of independent testing between single SVM and MDDLogo-clustered SVM models.....	27
Figure 19. Submitting the protein sequences in Fasta format interface for prediction.....	30
Figure 20. A case study of S-palmitoylation site prediction for CD9 antigen (CD9_HUMAN).....	31