

中國醫藥大學

醫務管理研究所碩士論文

論文編號: IHAS-292

序列組共同具有短序列樣式之詳盡搜尋演算法

Finding Motif by exhaustive algorithm



指導教授：莊 振 村 博士

研 究 生：楊 鎮 嘉 撰

中華民國九十三年六月

目 錄

第一章:緒論	
◆ 研究背景	..1
◆ 問題陳述	..2
◆ 研究問題	..4
◆ 研究目的	..6
第二章:文獻探討.....	7
第三章:研究設計與方法	
◆ 研究樣本.....	23
◆ 研究方法.....	25
◆ 理論步驟.....	26
◆ Xor 運算	28
◆ 比對技巧.....	29
◆ 研究方法實例.....	31
第四章:結論	37
第五章:未來發展	39
研究限制.....	41
參考文獻	42

表目錄

表 1.1	在不同序列長度下不同演算法的執行係數比較	4
表 1.2	不同演算法找尋不同類型 motif 的執行係數比較	5
表 3.1	Xor 運算	28
表 3.2	方法實例 pattern 紀錄表	31

圖目錄

圖 1.1 Challenge Problem 範例	4
圖 2.1 樹枝搜尋策略(1)	11
圖 2.2 樹枝搜尋策略(2)	11
圖 2.3 TEIRESIAS 演算法範例一	12
圖 2.4 TEIRESIAS 演算法範例二	13
圖 2.4 WINNOWER 演算法範例一	15
圖 2.5 WINNOWER 演算法範例二	16
圖 2.6 Consensus 演算法範例	17
圖 2.7 Random projection 步驟一	21
圖 2.8 Random projection 步驟二	21
圖 3.1 比對技巧一	29
圖 3.2 比對技巧二	30
圖 3.3 方法實例 S_5 與 S_2 比對後	32
圖 3.4 方法實例 S_5 與 $S_2 S_4$ 比對後	33
圖 3.5 方法實例 S_5 與 $S_2 S_4 S_3$ 比對後	34
圖 3.6 方法實例 S_5 與 $S_2 S_4 S_3 S_6$ 比對後	34
圖 3.7 方法實例 S_5 與 $S_2 S_4 S_3 S_6 S_1$ 比對	35
圖 3.8 方法實例最後結果	35

摘要

Motif Finding 這個問題引起許多人的興趣，因此也造就了許多搜尋演算法的產生。然而許多的演算法在面對Pevnzer 和 Sze所提出的Challenge Problem無法得到一個非常好的結果。而本研究主要目的是利用自行構思的詳盡式搜尋來解決Motif finding的問題，嘗試突破PROJECTION在搜尋 $(9,2)$ 、 $(11,3)$ 、 $(13,4)$ 、 $(15,5)$ 、 $(17,6)$ 等信號其結果不佳的問題，以及解決當序列長度(N)增加時所造成搜尋 motif 執行效率和準確度下降等問題。

一般人認為詳盡搜尋的方式雖然可以準確的找到 motif，但是當基因序列和搜尋的 motif 長度過長的時候，在運算時間容易呈指數的倍數成長，因此本研究透過一些輔助的技巧可以避免因長度的增加使得運算時間的增長，又能夠兼顧效率及精準的方式來找到最好的結果。可以預期的是此方式可以得到不錯的結果外，在方法上更可擴展到其他 motif 相關的研究領域上。

Abstract

Recently, motif finding became a very popular area in bioinformatics, thus more and more researches are interested in discover motif. However, many algorithms can not solve the Pevzner and Szec's challenge problem. It motivates my algorithm to purpose construct an exhaustive method to improve the motif finding performance in discover signals such as:

$(9,2)$, $(11,3)$, $(13,4)$, $(15,5)$, and $(17,6)$ and to solve the problem that accuracy will be descending while sequence length is increasing.

Although exhaustive search method could find motifs accuracy, it still needs to face the problem that computing time will be grown exponentially by length increasing of genomic sequence and motif. The research through provide assist skills not only to avoid the length of sequence increasing effect computing time, but also efficiency and accuracy to discover the optimal result. We could expect the research will bring well performance and apply in other bioinformatics domain related motif finding is expandable.

誌 謝

首先要感謝指導教授 莊振村博士引導我在踏入資訊的領域並且細心的指導與教悔，在這兩年的時間學習到做研究的精神、增進自我的邏輯與獨立思考的能力，此外，感謝吳禮字老師給予生物上的教學指導，如此才得以順利完成論文研究。另外，還要口試委員王旭正博士、周志賢博士於百忙之中審閱論文，並提供建議使得本論文更加完善。

在面臨問題的討論過程感謝梅芬在生物方面的知識的協助。最後，要感謝我的父母給我的支持與鼓勵，並且提供一個舒適安穩的求學環境，讓我能毫無牽掛的完成碩士學位。在此，謹將這份榮耀獻給我的家人。

研究生 楊鎮嘉 謹致

民國93年6月