

摘要

解碼基因為人類帶來了期望，因為可藉此了解遺傳及環境引起的人類進化、疾病和 DNA 內部的構造功能。人類基因組序列宣示將在 2005 年完成基因定序，為加速如此龐大的基因鹼基序列定序尋找基因，需仰賴電腦科技運算的基因預測程式，因此許多程式陸續地被開發。然而，基因預測工具相當多，目前這些程式都無法相當正確的尋找基因，還未能完全真正解決許多基因結構預測的問題。

本研究結合統計特徵的計算分析與候選 exon 區段計分演算兩種方式進行基因尋找。此方法的主要目的是定義 exon 序列及基因序列邊界，而其邊界是藉由比較基因序列區段的鹼基組成統計特徵，以及邊界的函數進行劃分。

本研究之基因預測的程式整合了電腦運算及統計技術，可在配合實驗室研究下幫助生物學家分析新定序的序列，並致力克服人類基因組複雜的結構，迅速搜尋基因，能以統計量化數據窺得序列全面的特徵，不需複雜的統計模式進行推估，仰賴訊號正確性，未來希望調查 exon 中鹼基出現頻率規則使得基因預測在實驗研究過程中化繁為簡，達成生物學家生物實驗研究先機，對於醫學研究影響深遠。

Abstract

Decoding gene bringing an expectation to understand how the heredity and environment affect human evolution, human disease and DNA structure and function. Human Genome Project publicizes that entire human gene sequencing will be finished in 2005; however the progress is uncertain so that many gene prediction programs are designed continuously to speed up gene finding in massive amount of genomic sequences. Although gene prediction programs are developed for the complex structure of human genomic sequences, but a accuracy program still be needed.

The research combined statistical characteristic analysis and exon region score algorithm to find gene region. The purpose of our method is to define exon region and gene region boundary. The boundary is generated by comparing base composition of gene region and scoring function of boundary.

The program integrates computer science and statistical analysis which tries to solve the problem of complex human genome structure could help biologists to find gene rapidly, the sequences are numeral to analysis global structure of genomic sequences without complex statistical model to depend on signal to estimate.

Future work is hopefully to help biologists to investigate exon frequency and discover exon composition rule to predict gene automatically.

誌謝

在論文完成的過程中，首先感謝指導教授莊振村老師在資訊方法上面臨困難時給予鼓勵，並且不斷的啟發自己去思考新的創意與靈感，影響我追求學問與研究的精神。此外，相當感謝口試委員王旭正博士及周志賢博士，百忙之中抽空審閱論文，給予指教，才能完成一份完善的研究論文。

本論文的研究需要相當多的分子生物知識，在充實分子生物的知識時，感謝吳禮宇老師及林妍如老師以其專業知識為不懂的地方耐心解惑。

最後在完成論文的討論過程中，感謝鎮嘉對於方法上的建議與意見，並且對於完成論文過程中不斷給予支持與鼓勵，也感謝在我完成論文的期間所有關心過我的朋友與同學，最後感謝我的父母提供一個舒適安穩的環境，使我無後顧之憂的完成碩士學位，因此將這份榮耀獻給我的家人。

研究生 黃梅芬 謹致

中華民國九十三年六月

目錄

第一章 緒論

- 一、研究背景1
- 二、問題陳述2
- 三、研究問題6
- 四、研究目的8

第二章 文獻探討

- 一、基因預測程式的探討9
- 二、外部排比法及內部排比法的比較24

第三章 研究設計與方法

- 一、研究假設28
- 二、研究模型29
- 三、資料來源30
- 四、測量方法30

第四章 討論43

參考文獻45

圖目錄

圖 1-1 exon 和 intron 有不一樣的 GC 比例	6
圖 1-2 人類基因表達的遺傳訊息.....	7
圖 2-1 HMM 中的隱藏狀態	12
圖 2-2 簡單 HMM 模式	13
圖 2-3 GENSCAN 的一般基因組序列的結構模式	15
圖 2-4 Grail II 編碼區域的認可模式圖(Coding recognition module; CRM)	19
圖 3-1 根據 exon 觀察次數作為計分依據的範例.....	38
圖 3-2 函數中加入觀察 exon 模糊區段的容忍值 Q 的範例.....	38
圖 3-3 容忍值 Q 為 1 時，會出現 2 段的候選 exon 的範例.....	38
圖 3-4 第 2 段 intron 為位置 54-58 的 gtaat.....	41
圖 3-6 函數中加入觀察 intron 模糊區段的範例.....	42
圖 3-7 調整後第 2 段 intron 位置應為 54-62 的 gtaatcaag.....	42

表目錄

表 1-1 生物基因體之定序與基因密度.....	3
表 2-1 內部法各種基因預測程式精準度比較.....	26
表 3-1 Genome table 資料型態.....	32
表 3-2 研究方法範例之 Genome table.....	33
表 3-3 exon table 資料型態.....	34
表 3-4 研究方法範例之 Exon table.....	35
表 3-5 基因組與 exon 比較表，表格名稱 Genomexon table.	35
表 3-6 研究方法範例之 Genomexon table	36
表 3-7 整理 exon 特徵區段計分表的資料型態，表格名稱:exonscore table.....	39
表 3-8 研究方法範例 gtgcaattctcccaaac 之 exonscore table.....	40
表 3-9 容忍值為 5 時，表格合併結果	41