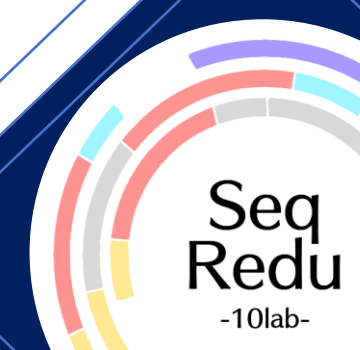


高效能生物資料去冗餘演算法

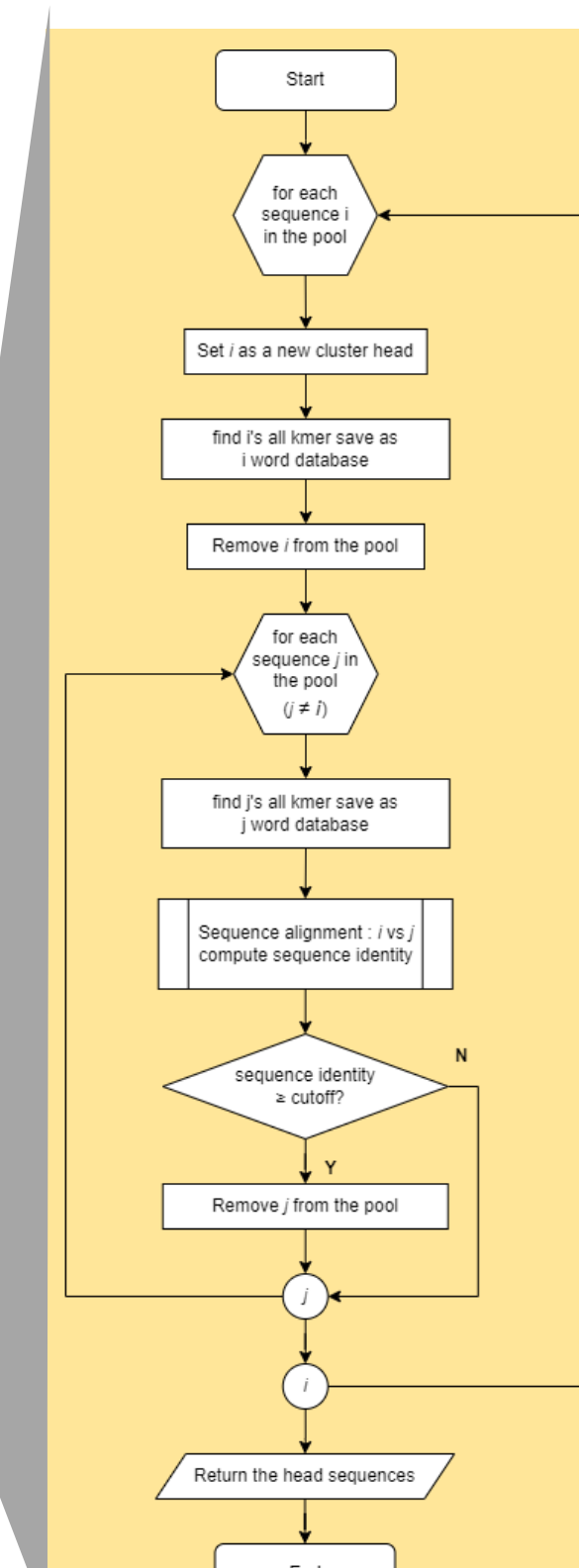
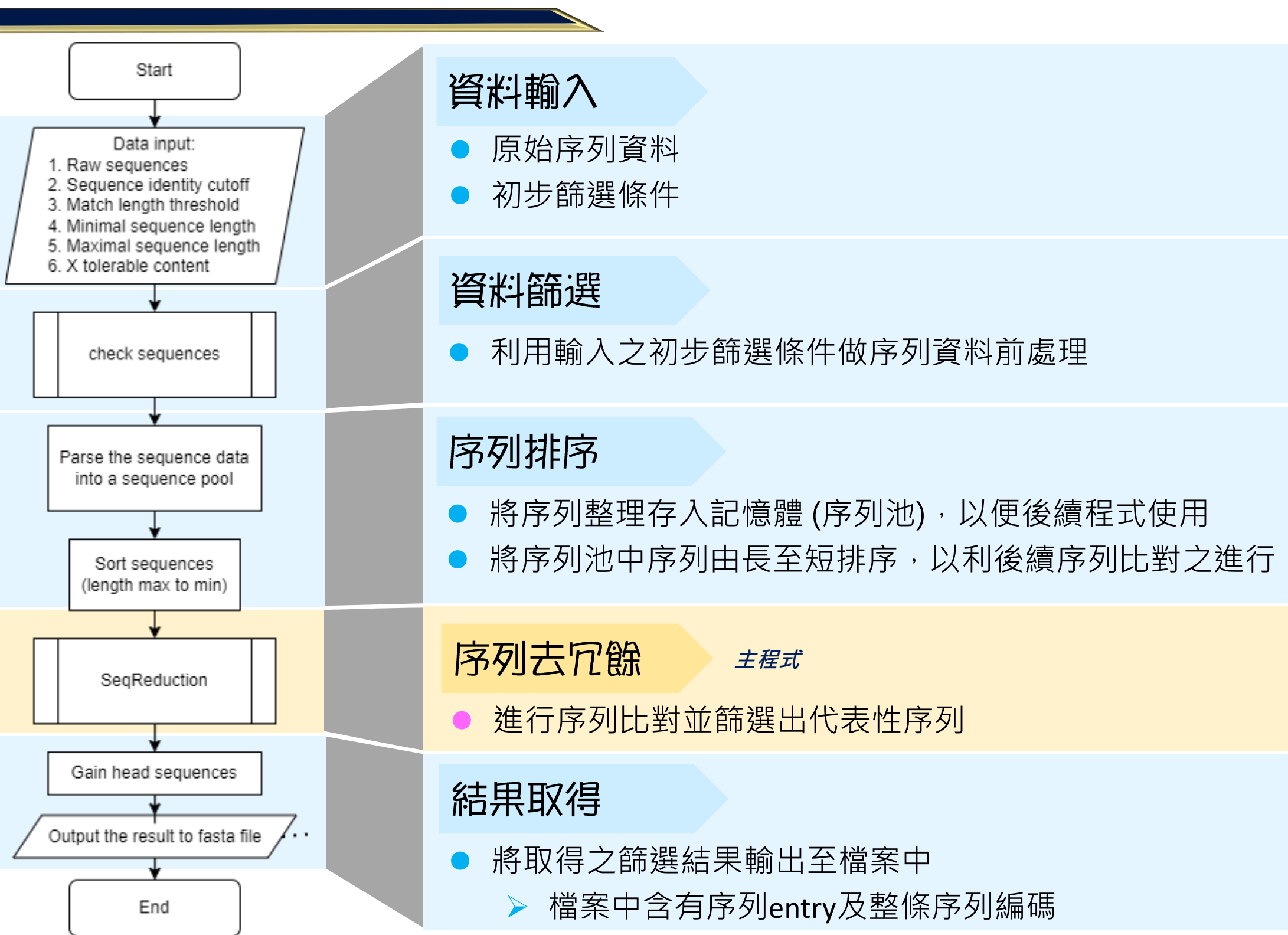
第十一組 高效能生物資料去冗餘演算法開發小組



想法理念

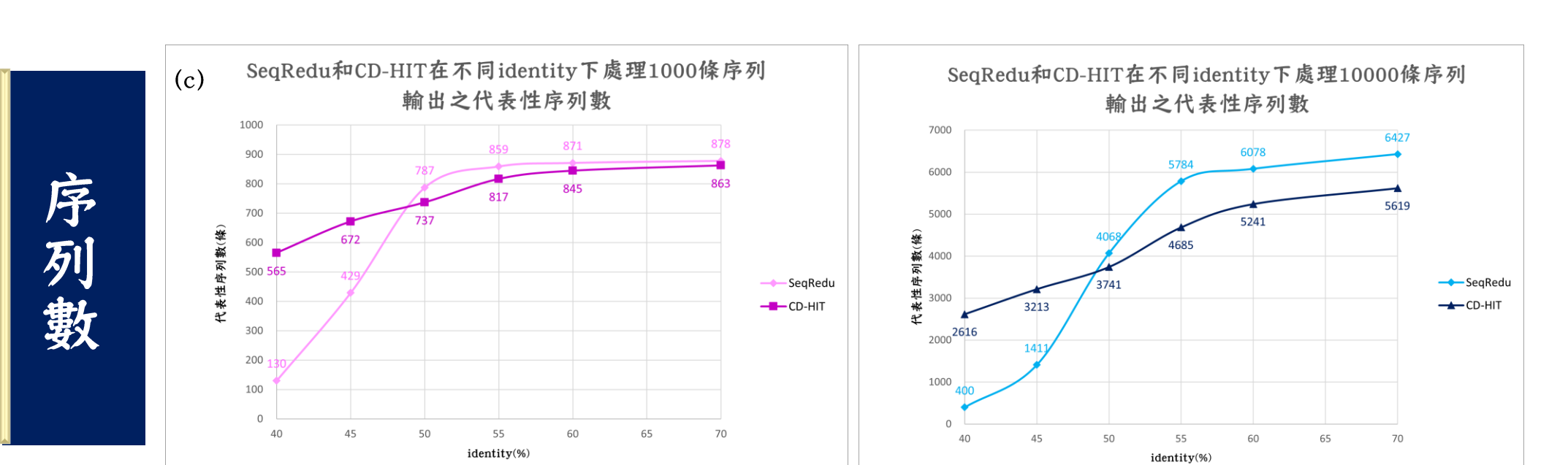
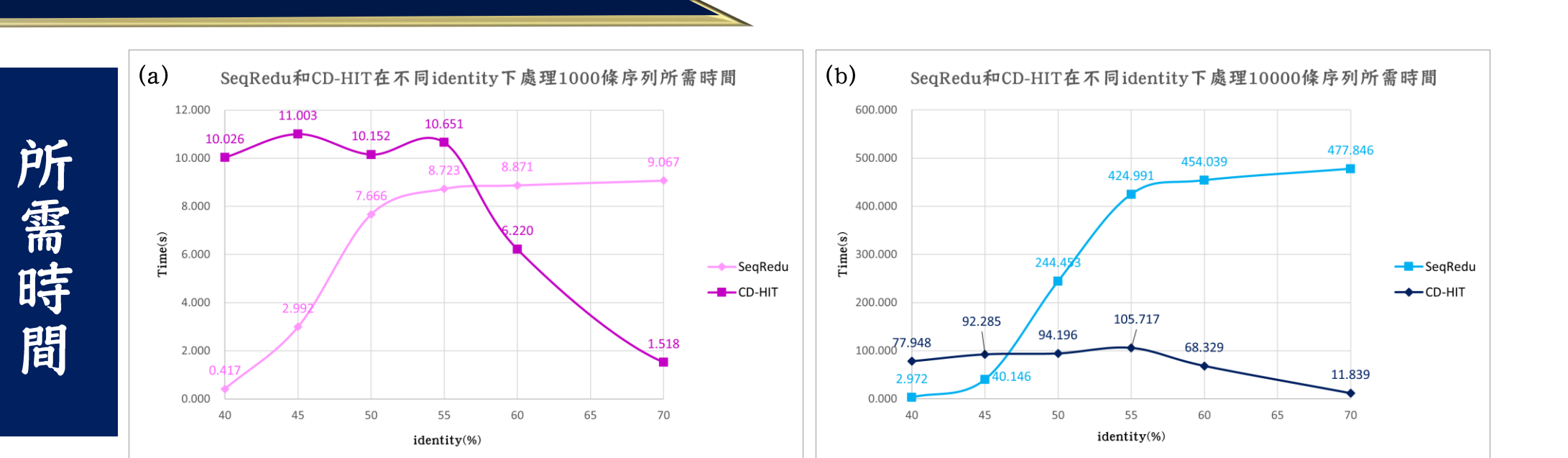
我們相信，臺灣作為資訊科技巨頭，有實力帶領全球開創數位生物科技新紀元。我們企圖拋磚引玉，以提升大數據分析效率為切入點，開發高效能生物序列資料去冗餘演算法，盼加速全球生物科技發展，並鼓舞更多資訊團隊投入生物科技研發，提升臺灣的國際地位。大數據時代來臨，處理海量數據極耗資源與時間，但許多生物資料具高度相似性，全數分析並無必要。當前諸多生物資料去冗餘演算法仍有缺點，速度快者不夠精確，精確者運算速度卻難以跟上資料增加腳步。我們發現，幾乎所有序列去冗餘演算法都將去冗餘及分群共同封裝，但很多研究只需去冗餘之資料而非分群結果。因此，我們欲開發一演算法，把傳統同時進行去冗餘和分群的操作做兩階段拆分，並於去冗餘階段用 Golang 程式語言實現高度平行運算，來提升資料生產效率。期望在海量數據充斥的現今，能以臺灣強大資訊實力，提供全球一個快速獲取高密度原始資料並銳減資料總量之有效解決方案，推進各生物科技領域研發，共創新紀元。

演算流程



- **第一層迴圈：**
取出序列池中最開頭的序列直接當作代表性序列，計算並儲存其文字kmer(為某序列中某長度的序列短片段。如一序列：AGDEE，其長度為三的kmer有AGD、GDE、DEE)。
- **第二層迴圈：**
運用分散運算一次進行多執行緒計算，透過和其他序列進行序列比對後得出之相似度進行篩選，如相似度大於使用者設定值，此序列將被視為和代表性序列有一定相似度，而從序列池中被移除。
- **所有執行緒都完成計算：**
進入下回合迴圈，重複直至序列池中所有序列都已比對完成。過程中如遇到已計算過文字kmer的序列，將會從緩衝區中取出先前計算之結果，直接進行比對。
- **Sequence alignment副程式：**
穿繞兩序列kmer集合資料，如有相同文字片段，則會透過已儲存在其中的文字片段位置，往下進行BLAST(Basic Local Alignment Search Tool)序列比對算法的計算，切除頭尾分數較低之配對後，再計算兩序列相似度

研究結果



隨著處理序列數量上升(a至b圖)，在低identity時，SeqRedu的計算速度可以較CD-HIT快上許多，尤其當資料量越多時，耗費時間差距會越顯著；但在高identity時，SeqRedu所耗費之時間會急遽上升，尤其在處理序列數為10000條(b)時，CD-HIT計算70% identity只需約11秒，而SeqRedu所耗費的時間卻高達8分鐘左右。因此，如何優化SeqRedu，提升其在高identity時的計算效能，將是我們在未來開發過程中，需要彼此充分思考討論並解決的一大關卡。

在高identity時，兩演算法輸出之代表性序列數會相近，SeqRedu輸出之序列數會稍多於CD-HIT，但當identity降低(由右至左)時，兩個演算法在不同資料量大下，皆會在約50% identity處發生交叉，當identity低於50%，SeqRedu所篩選出之序列數相較CD-HIT會有明顯的降低。因為SeqRedu在計算過程中不會開gap(gapless alignment)，在程式停止後切除頭尾分數較低之配對，並直接計算兩序列相似度，最後再利用輸出之相似度進行序列去冗餘暨分群。因此，SeqRedu較不容易算出高identity的結果。

網站成果展示

網站QR code

參與人員

SeqRedu 網頁大綱

- 首頁 (Home)
 - 背景簡介
 - 演算法簡介
 - 從生物序列到蛋白質
 - 去冗餘&生物序列分析
- 介紹 (Introduction)
 - 演算法詳細說明介紹
 - 10lab實驗室網頁連線
 - 作者介紹
- 下載 (Download)
 - 演算法檔案下載
 - 舊版SeqRedu執行檔
 - 新版SeqRedu執行檔
- 參考資料 (Doc & Reference)
 - 相關演算法網頁
 - 相關內容參考文章
- 聯絡我們 (Contact us)
 - 相關建議填寫區
 - 實驗室相關聯絡資訊

Sequence Reduction
高效能生物資料序列去冗餘演算法

這是由鈞鈞、千珊、詩玟、啓宏、曉典，在羅惟正教授及梁美智教授的指導下，做出的研究成果。

我們的目標是：開發新時代的生物資訊演算法，提升台灣在生物科技的國際地位

<p>Associate Professor</p> <p>Wei-Cheng Lo (羅惟正)</p> <p>WadeLo@nycu.edu.tw</p>	<p>Associate Professor</p> <p>Mei-Chih Liang (梁美智)</p> <p>mcliang@nycu.edu.tw</p>	<p>Undergraduate</p> <p>Chun-Yi Yang (楊鈞詒)</p> <p>neoccat.bt09@nycu.edu.tw</p>	<p>Undergraduate</p> <p>Chien-Shan Lin (林千珊)</p> <p>lcs0909.bt09@nycu.edu.tw</p>	<p>Undergraduate</p> <p>Shih-Wen Weng (翁詩玟)</p> <p>zaki.bt09@nycu.edu.tw</p>	<p>Undergraduate</p> <p>Yang-Tien Chou (周曉典)</p> <p>csowo.cs09@nctu.edu.tw</p>	<p>Undergraduate</p> <p>Qi-Hong Su (蘇啓宏)</p> <p>mequmi.bt09@nycu.edu.tw</p>
--	---	--	--	--	--	---