生命情報学(2)

配列解析基礎

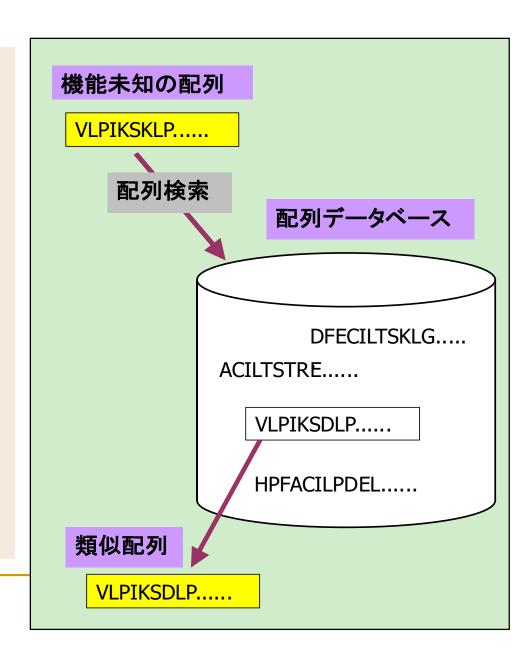
阿久津 達也

京都大学 化学研究所 バイオインフォマティクスセンター

配列アラインメントとは?

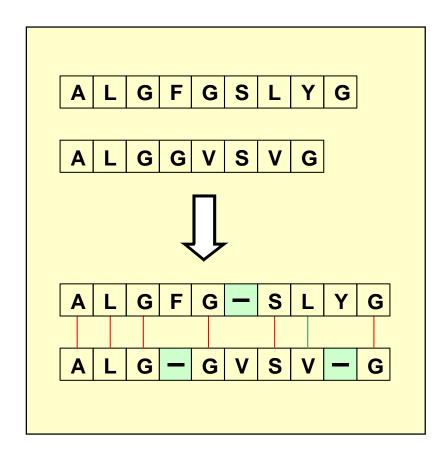
配列検索

- バイオインフォマティク スにおける基本原理
 - □ 配列が似ていれば機能 も似ている
 - □ ただし、例外はある
- 配列検索の利用法
 - □ 実験を行い機能未知の配列 が見つかった
 - □ データベース中で類似の配列 を検索
 - 機能既知の類似の配列が見つかれば、その配列と似た機能を持つと推定



配列アラインメント

- バイオインフォマティクスの 最重要技術の一つ
- ・ 2個もしくは3個以上の配列の類似性の判定に利用
- 文字間の最適な対応関係 を求める(最適化問題)
- 配列長を同じにするように 、ギャップ記号(挿入、欠失 に対応)を挿入



2個の配列に対するアラインメント: ペアワイズ・アラインメント 3個以上の配列に対するアラインメント: マルチプル・アラインメント

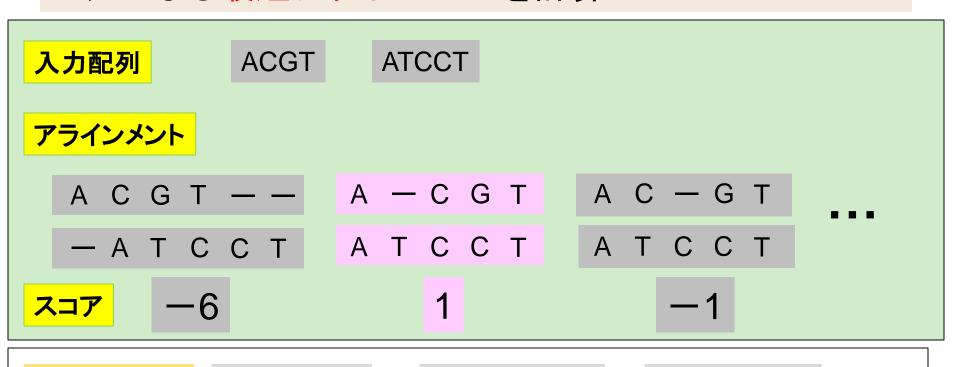
ペアワイズ・アラインメント

ペアワイズ・アラインメント

同じ文字: 1

スコアの定義

- 2本の配列に対するアラインメント
 - 大域アラインメント: 配列全体にわたるアラインメント
- 列ごとにスコアが定義され、各列のスコアの和が最大となる最適アラインメントを計算

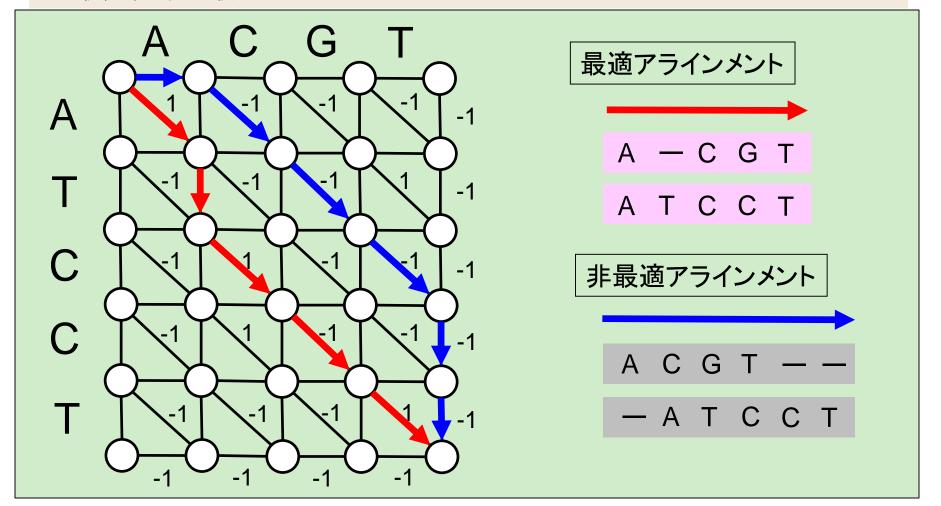


違う文字: -1

ギャップ: -1

大域アラインメントと格子状グラフ

- 入力文字列から格子状グラフを構成 (縦横の辺の重みはすべてギャップスコア、斜めの辺は対応する文字間のスコア)
- アラインメントと左上から右下へのパスが一対一対応
- 最長経路=最適アラインメント



動的計画法による最適アラインメントの計算

- アラインメントの個数:指数関数のオーダー
- 動的計画法を用いれば *O(mn)* 時間
 - D[i,j] は始点(0,0)から(i,j)までの最適パスのスコア
- アラインメントの復元(トレースバック)
 - *D*[*m*,*n*] から再帰式で=となっている頂点を逆にたどる

$$D[i,0] \leftarrow i \times (-d) \qquad i = 0, \dots, m$$

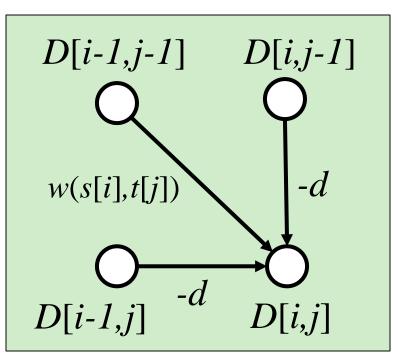
$$D[0,j] \leftarrow j \times (-d) \qquad j = 0, \dots, n$$

$$D[i-1,j] - d$$

$$D[i,j-1] - d$$

$$D[i-1,j-1] + w(s[i],t[j])$$

s,t: 入力配列 s[i]: 配列 s の i 番目の文字 m=|s|, n=|t| w(x,y): 文字 x,y 間のスコア -d: ギャップ記号のスコア



スコア行列

М

S

7/

- 残基間(アミノ酸文字間)の類似性を表す行列
 - PAM250, BLOSUM45 など

BLOSUM50 スコア行列 (置換行列)の一部分

局所アラインメント

局所アラインメント

- 配列の一部のみ共通部分があることが多い
 - ⇒共通部分のみのアラインメント
 - □ 例えば、AATGCATT と GATCG の場合、

ATG C

AT-C

というアラインメントを計算

- ■問題の定義
 - □ 入力: 2個の配列 *s*, *t* スコア関数 *w*(*x*, *y*)
 - □ 出力: S_{opt}(s[h...k],t[h'...k']) が最大となる部分文字列の組(s[h...k],t[h'...k'])に対する最適アラインメント
- 大域アラインメントを繰り返すとO(m³n³)時間
 - ⇒Smith-Watermanアルゴリズムなら*O(mn)*時間

局所アラインメントに対する動的計画法

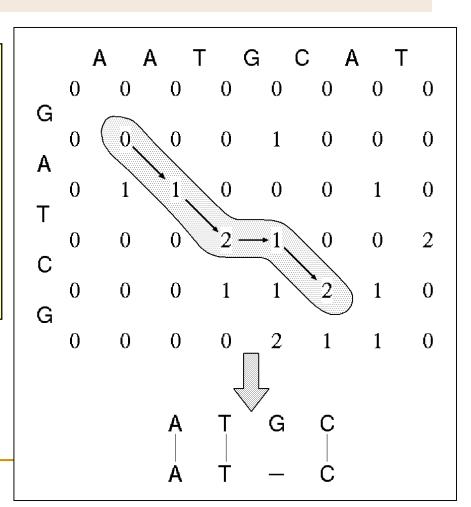
大域アラインメントに対する動的計画法を少し修正するだけでOK

$$D[i,0] \leftarrow 0 \qquad i = 0, ..., m$$

$$D[0,j] \leftarrow 0 \qquad j = 0, ..., n$$

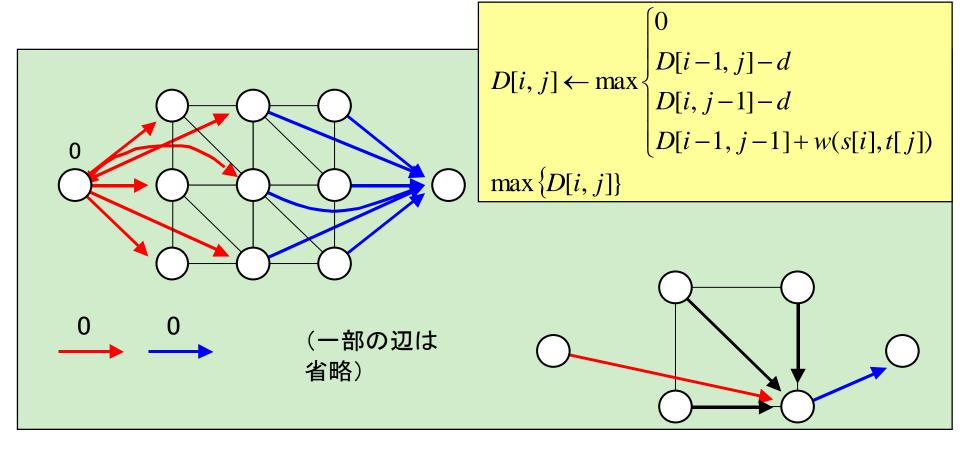
$$D[i,j] \leftarrow \max \begin{cases} 0 \\ D[i-1,j] - d \\ D[i,j-1] - d \\ D[i-1,j-1] + w(s[i],t[j]) \end{cases}$$

$$\max \{D[i,j]\}$$



局所アラインメント・アルゴリズムの正当性

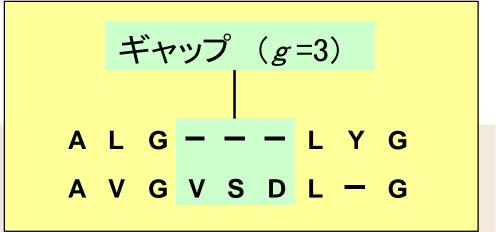
- 証明のアイデア
 - □ 始点と終点を表す2個の頂点を格子状グラフに追加
 - □ 始点から終点へのパスと局所アラインメントが1対1対応



ギャップコスト

ギャップペナルティ

- 線形コスト -gd
 - g: ギャップ長
 - □ d: ギャップペナルティ
 - □ この図の例では、コスト= -3d
- アフィンギャップコスト -d e(g-1)
 - □ d: ギャップ開始ペナルティ
 - □ e: ギャップ伸張ペナルティ
 - □ この図の例では、コスト= -d 2e
 - □ よく利用されるペナルティ (d,e)=(12,2),(11,1)



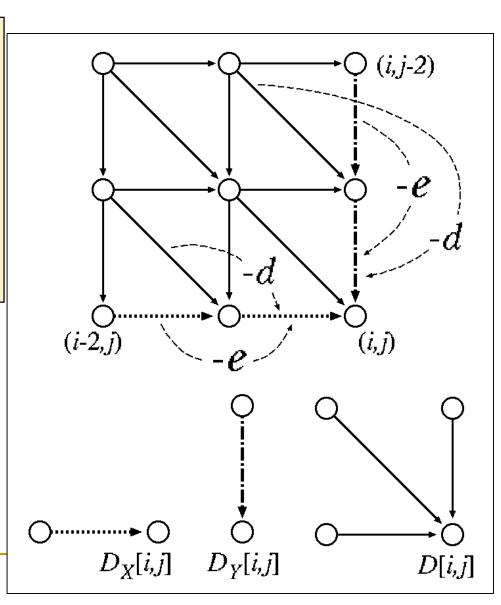
アフィンギャップコストによるアラインメント

$$D_{X}[i, j] = \max \begin{cases} D[i-1, j] - d \\ D_{X}[i-1, j] - e \end{cases}$$

$$D_{Y}[i, j] = \max \begin{cases} D[i, j-1] - d \\ D_{Y}[i, j-1] - e \end{cases}$$

$$D[i, j] = \max \begin{cases} D_{X}[i, j] \\ D_{Y}[i, j] \\ D[i-1, j-1] + w(s[i], t[j]) \end{cases}$$

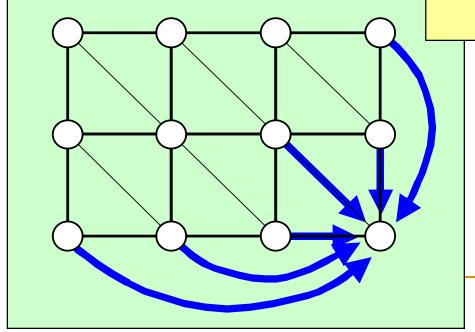
- 三種類の行列を用いる動的 計画法によりO(mn)時間
- Smith-Watermanアルゴリズムとの組み合わせが広く利用されている
- ⇒ Smith-Waterman-Gotoh アルゴリズム



任意ギャップコストによるアラインメント

■ 動的計画法(下式)により、 $O(n^3)$ 時間 (ただし、m=O(n)とする)

$$D[i, j] \leftarrow \max \begin{cases} D[i-1, j-1] + w(s[i], t[j]) \\ \max_{k=0, \dots, i-1} D[k, j] + \lambda(i-k) \\ \max_{k=0, \dots, j-1} D[i, k] + \lambda(j-k) \end{cases}$$



配列検索の実用的アルゴリズム

配列検索の実用プログラム(1)

- O(mn): m は数百だが、n は数GBにもなる
 - ⇒実用的アルゴリズムの開発

FASTA: 短い配列(アミノ酸の場合、1,2文字、DNAの場合、4-6文字)の完全一致をもとに対角線を検索し、さらにそれを両側に伸長し、最後にDPを利用。

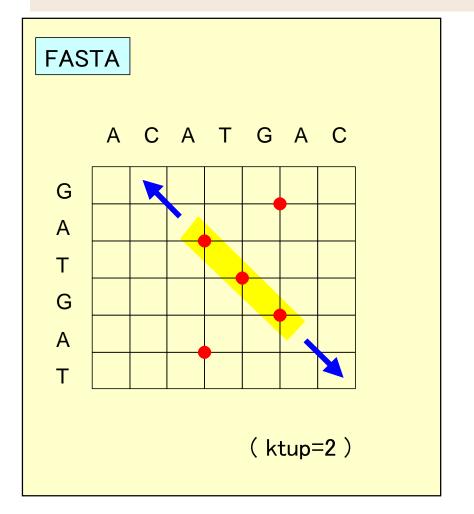
• BLAST: 固定長(アミノ酸では3, DNAでは11)の全ての類似単語のリストを生成し、ある閾値以上の単語ペアを探し、それをもとに両側に伸長させる。ギャップは入らない。伸長の際に統計的有意性を利用。

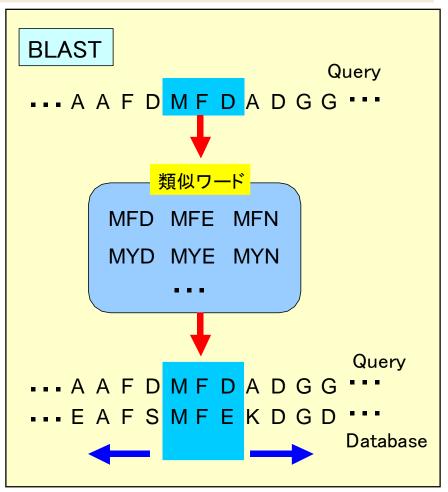
配列検索の実用プログラム (1)

- データベース検索に O(mn)時間: m は数百だが、n は数GBにもなる
 - ⇒実用的アルゴリズムの開発
- FASTA: 短い配列(アミノ酸の場合、1,2文字、DNAの場合、4-6文字)の完全一致をもとに対角線を検索し、さらにそれを両側に伸長し、最後にDPを利用。
- BLAST: 固定長(アミノ酸では3, DNAでは11)の全ての類似単語のリストを生成し、ある閾値以上の単語ペアを探し、それをもとに両側に伸長させる。ギャップは入らない。伸長の際に統計的有意性を利用。

配列検索の実用プログラム (2)

- FASTA: 短い配列(アミノ酸の場合、1,2文字、DNAの場合、4-6文字)の完全一致をもとに対角線を検索し、さらにそれを両側に伸長し、最後にDPを利用。
- BLAST: 固定長(アミノ酸では3, DNAでは11)の全ての類似単語のリストを生成し、ある閾値以上の単語ペアを探し、それをもとに両側に伸長。





配列検索の実用プログラム (3)

- SSEARCH: 局所アラインメント(Smith-Waterman-Gotohアルゴリズム)をそのまま実行
- PSI-BLAST: ギャップを扱えるように拡張したBLAST を繰り返し実行。「BLASTで見つかった配列からプロファイルを作り、それをもとに検索」という作業を繰り返す。
- PatternHunter: 穴あきシードを用いる(連続した文字ではなく飛び飛びの文字の完全一致をもとに検索)

マルチプル・アラインメント

マルチプル・アラインメント: 意味

- 3本以上の配列が与えられた時、全ての配列の長さが 同じになるようにギャップを挿入
- 進化的、構造的に相同な残基(塩基)ができるだけ同じ カラムに並ぶようにする
- 通常はスコアを用いて、最適化問題として定式化
- 理想的なアライメント
 - □ 同一残基から派生した残基が同一カラムに並ぶ
 - □ 構造的に重なり合う残基が同一カラムに並ぶ
 - ⇒構造的に重なり合わない場所を無理に重ね合わせるのは、あ まり意味がない

マルチプル・アラインメント: 定式化

3本以上の配列が与えられた時、長さが同じで、かつ、スコアが最適となるように各配列にギャップを挿入したもの

HBA_HUMAN VGAHAGEY
HBB_HUMAN VNVDEV
MYG_PHYCA VEADVAGH
GLB5_PETMA VYSTYETA
LGB2_LUPLU FNANIPKH
GLB1_GLYDI IAGADNGAGV



HBA_HUMAN V G A - - H A G E Y HBB_HUMAN V - - - N V D E V MYG_PHYCA V E A - - D V A G H GLB5_PETMA V Y S - - T Y E T A LGB2_LUPLU F N A - - N I P K H GLB1_GLYDI I A G A D N G A G V

- \bullet スコアづけ (全体スコアは基本的に各列のスコアの和: $\sum S(m_i)$)
 - □ 最小エントロピースコア
 - $S(m_i) = -\sum c_{ia} \log p_{ia}$ ($c_{ia} = i$ 列におけるaの出現回数, $p_{ia} = i$ 列におけるaの生起確率)
 - □ SPスコア(Sum-of-Pairs)
 - $S(m_i) = \sum_{k < l} w(m_k[i], m_l[i])$

 $(m_k[i]=$ アラインメント後のi列, k行目の文字)

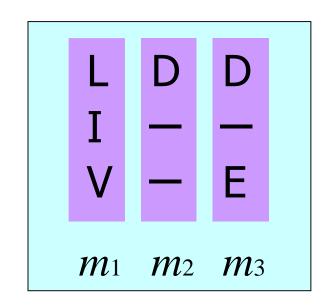
SP (Sum of Pairs) スコア

- $S(m_i) = \sum_{k < l} w(m_i^k, m_i^l)$ $m_i^k = i \mathcal{M}, k \mathcal{T} = \mathcal{O}$
- ■問題点
 - □ 確率的な正当性が無い
 - 。同一カラムに a,b,c が並んだ場合、 $\log(p_{abc}/q_aq_bq_c)$ とすべきだが、SPスコアでは $\log(p_{ab}/q_aq_b) + \log(p_{bc}/q_bq_c) + \log(p_{ac}/q_aq_c)$

$$S(m_1) = w(L, I) + w(L, V) + w(I, V)$$

$$S(m_2) = w(D, -) + w(D, -) + w(-, -)$$

$$S(m_3) = w(D, -) + w(D, E) + w(-, E)$$



多次元DPによるマルチプル・アラインメント

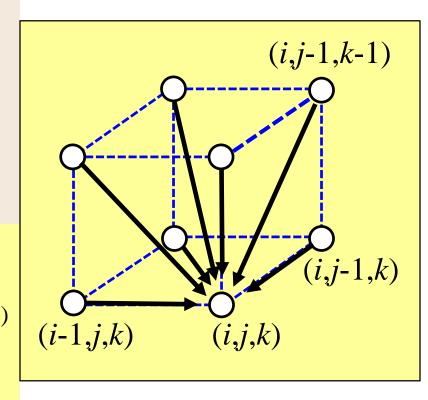
N個の配列に対するマルチプル・アラインメント

N次元DPにより $O(2^N n^N)$ 時間 (各配列の長さはO(n)を仮定)

■ 例:N=3

 $D[i, j, k] \leftarrow \max \begin{cases} w(s_1[i], s_2[j], s_3[k]) \\ D[i, j-1, k-1] + w(-, s_2[j], s_3[k]) \\ D[i-1, j, k-1] + w(s_1[i], -, s_3[k]) \\ D[i-1, j-1, k] + w(s_1[i], s_2[j], -) \\ D[i, j, k-1] + w(-, -, s_3[k]) \\ D[i, j-1, k] + w(-, s_2[j], -) \\ D[i-1, j, k] + w(s_1[i], -, -) \end{cases}$

D[i-1, j-1, k-1] +



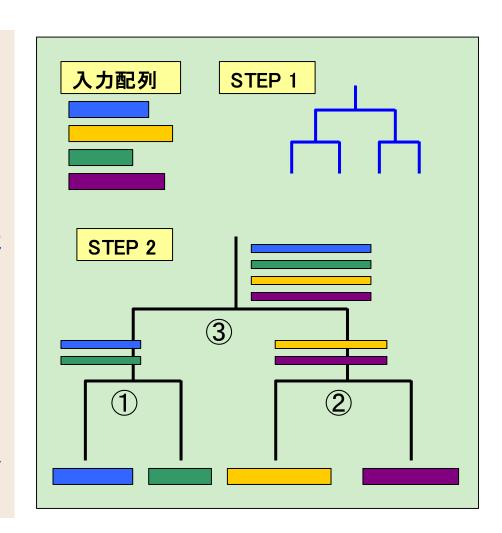
■ 一般の N に対しては NP困難

マルチプル・アラインメントの実用的計算手法

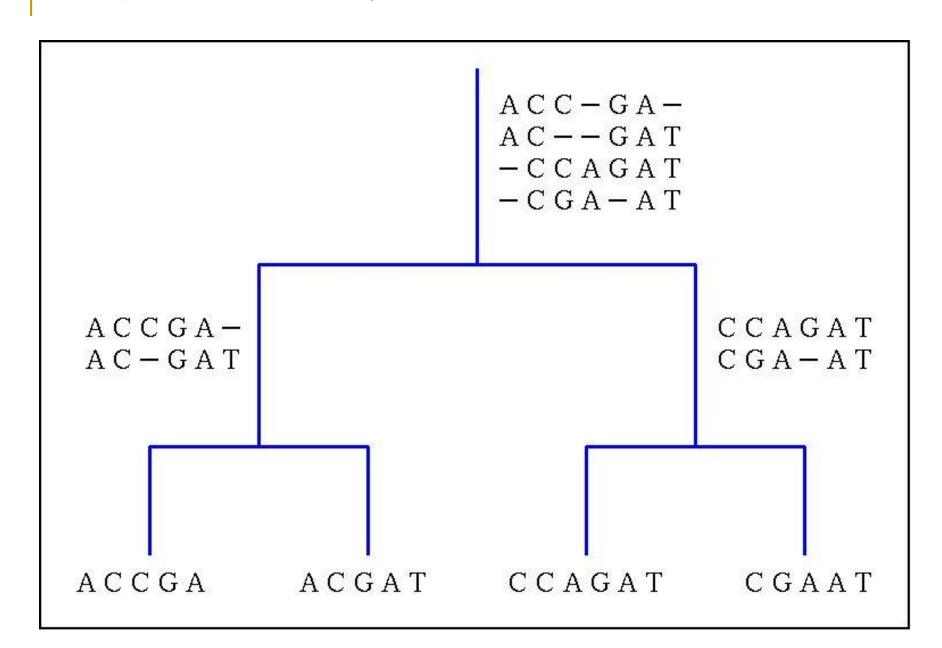
- プログレッシブ・アラインメント
 - □ CLUSTAL-W(広く利用されているソフト)などで採用
 - □ 逐次改善法との組み合わせが、より有効
- 逐次改善法
- シミュレーテッドアニーリング
- 遺伝的アルゴリズム
- HMMによるアラインメント
- ■分枝限定法
 - □ 10配列程度なら最適解が計算可能な場合がある

実用的マルチプル・アラインメント法

- ヒューリスティックアルゴリズムの 開発
 - □ *N*次元DPは(*N*=4ですら) 非実用的
 - □ 一般にはNP困難
- プログレッシブアラインメント
 - 1. 近隣結合法などを用いて 案内木 を作る
 - 類似度が高い節点から低い節点へという順番で、配列対配列、配列対プロファイル、プロファイル対プロファイル対プロファイルのアラインメントを順次計算
- 逐次改善法
 - 「配列を一本取り除いては、アライ ンメントしなおす」を繰り返す

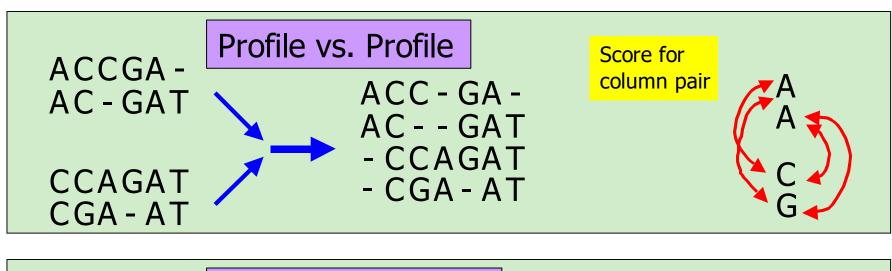


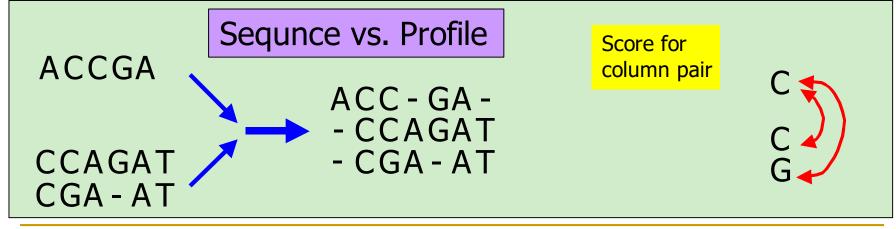
プログレッシブ・アラインメント



プロファイループロファイル・アラインメント

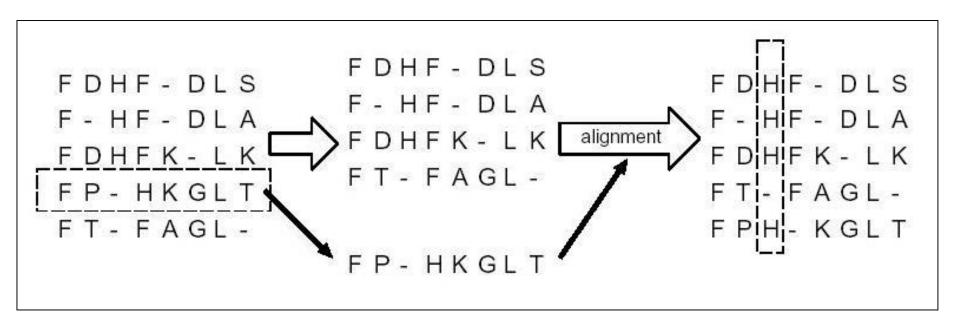
■ 各列を1文字のように扱うことにより、DPにより計算





逐次改善法

「配列を一本取り除いては、アラインメントしなおす」 を繰り返す



まとめ

- ペアワイズ・アラインメント
 - □ 大域アラインメント
 - 最長パス問題に変換し、動的計画法を適用
 - *O*(*mn*)時間
 - □ 局所アラインメント
 - 類似部分のみのアラインメントを計算
 - アフィンギャップコストを用いたアラインメント
 - ■「一度ギャップが入ると連続して入りやすい」を反映
- マルチプル・アラインメント
 - □多次元DP
 - N本の配列のアラインメント $\Rightarrow N$ 次元動的計画法
 - 最適解が計算可能だが O(2^Nn^N) 時間かかり非実用的
 - □ プログレッシブ・アラインメント
 - 最適性の保証はないが実用的