

ICOT Technical Memorandum: TM-1212 他

---

TM-1212 他

情報処理学会  
第45回全国大会論文集

October, 1992

© 1992, ICOT

**ICOT**

Mita Kokusai Bldg. 21F  
4-28 Mita 1-Chome  
Minato-ku Tokyo 108 Japan

(03)3456-3191~5  
Telex ICOT J32964

---

**Institute for New Generation Computer Technology**

TM1212	知識を用いた蛋白質配列解析システムの試み	廣澤 誠、星田 昌紀、石川 幹人
TM1213	PIM/m のメンテナンス系アーキテクチャと CSP	高橋 勝己、武田 保孝、村沢 靖、小森 隆三（三菱）、杉野 荣二（北陸先端科学技術大学院大学）
TM1214	並列シミュレーションアーリングとタンパク質配列解析	戸谷 智之、石川 幹人、星田 昌紀、荒木 均（松下）
TM1215	並列反復改善法によるタンパク質配列のアライメント	星田 昌紀、石川 幹人、廣澤 誠、戸谷 智之
TM1217	階層的な文章構造表現に基づく接続表現の生成と省略の処理	池田 光生
TM1220	法的推論システム HELIC-II (3) 一判例を用いた類似検索と判断生成	前田 茂、小野 昌之、新田 克己
TM1223	並列データベース管理システム Kappa-P	河村 元夫、横田 一正、佐藤 裕幸、永沼 和智、中嶋 かおり、椿野 宣行、川村 達（三菱）、澤部 直太、西山 聰、富樫 泰子（三菱総研）、坂下 之子、合田 光宏 (ア・ティフィシャル・インテリジェンス)
TM1224	演繹オブジェクト指向データベースの分子生物学への摘要	田中 秀俊
TM1228	法的推論システム HELIC-II (2) 類似事例検索の改良と評価	小野 昌之、前田 茂、新田 克己

## 知識を用いた蛋白質配列解析システムの試み

廣沢 誠、星田 昌紀、石川 幹人

(財) 新世代コンピュータ技術開発機構

### 1はじめに

蛋白質の相間性解析の技術であるマルチブル・アライメントは、蛋白質の機能、構造子測、生物種の進化系統樹の作成の際に欠かせない技術である。従来は、複数本の蛋白質配列のマルチブル・アライメントは、生物学者が経験と勘を頼りに行ってきました。しかし、蛋白質配列の決定技術が著しく進歩したために、アライメントするべき蛋白質配列のみではなく、アライメントの回数も増えてきました。このため、計算機を用いたマルチブル・アライメントが導入されつつある。

現在まで多くのマルチブル・アライメントのアルゴリズムが開発されてきている。これらは、アライメントに対して定義された評価値を最適化することを目的とするものである。配列の本数が少ない時(2 or 3)には上記の計算機的に最適なアライメントを求めることができる[Needleman and Wunsch 1978; Murata 1985]。それ以上の本数の場合にも理論的には最適なアライメントを求めることが可能であるが、計算量が膨大であるので、実際には準最適なアライメントを求めるアルゴリズムが用いられる[Barton 1990; Berger and Manson 1991]。

しかしながら、上記のアルゴリズムが導き出すアライメントは、生物学的に意味のあるアライメントでは必ずしもない。我々は、原点に戻り、生物学的に意味のあるアライメントとは何かを考えた。そして、アライメントの専門家にインタビューし、どのようなアライメントを良いアライメントであると見なしているかを把握し、また、彼らがアライメントを行う時に意識、無意識で用いているルール、知識を抽出した。そして、これらを解析した結果を反映したマルチブル・アライメントシステムを試作した。

以下、計算機的に最適なアライメントが、生物学的に意味のあるアライメントではない例を示す。そして、我々のマルチブル・アライメントシステムを紹介し、このアライメントシステムが生物学的に意味のあるアライメントを作成することを示す。なお、このシステムの詳しい内容については[Hiroswawa 1992]を参照していただきたい。

### 2計算機的に最適なアライメントの問題点

この章では、計算機的に最適なアライメントが、生物学的に意味のあるアライメントではない例を示す。しかしながら、生物学的に意味のあるアライメントの定義をアライインされる配列に依存せずに行うこととは困難である。ここでは、レトロウイルスというウイルスに特有である endonuclease という蛋白質のマルチブル・アライメントを例題にとる。なぜなら、この蛋白質の生物学的に意味のあるアライメントは、過去の研究により明らかであるからである。図1に6種類のレトロウイルスが持つ endonuclease の蛋白質配列をアライメントしたものを見ると、生物学的に意味のあるアライメントではない例が示されている。

アライメントのひとつである。

図1のアライメントが生物学的に意味のあるアライメントである理由は、図にも示されているように、全ての配列において共通するアミノ酸として、左側に二つの“H”、右側に二つの“C”を捕らえているからである。分子生物学では、全ての配列に共通なアミノ酸パターンなどをモチーフと呼ぶが、この図のモチーフは、Zinc Finger のモチーフと呼ばれており、生物学的に重要な機能を持つとされている。以降、例の蛋白質配列のアライメントでは、Zinc Finger のモチーフを持つものを生物学的に意味のあるアライメントと定義する。

```
17.6 : -----LLD-F-----KQDLVQDQTK-LF-GCT-TY-FPFLQLQGQIINCSL-AKT-DRS-T-TMPXTT
M-MULV : -----LLD-FL-----RQ-LTE-SPLRS-LVLLERSSPPMMRQMLT-KRLLTOMAGD-YAA-SQS-----+---+---+
RSV : LLDALL-TTP-VLG-EQPLRS-FTTNGGATL-T-LD-----CATTGTA-STILASCGAAG-QGKQPKC-----+
HIV : VASQGATGQ-FPLA-EAQQELT-ALGGQPLA-Q-HA-----CH13H0H-APVYQTCPC-BSA-FALAEQ-VI-
PRV : LSD-PFL-EATQAT-GLQE-AATL-B-LL-----TLLTAAQH-ADTVKAQEQCV-TPPFLG-B-VI-
SIV : LLL-T-NLA-SAQESKA-LQGQHAAAL-A-PG-----TATTAQH-ADTVLGCPFPOGSA-FQI-B-VI-
```

図1：生物学的に意味のあるアライメントの例

前に述べたように、計算機的に最適なアライメントを求められるのは、3本の配列のアライメントまでである。これは、Dynamic Programming という手法で行われる[Needleman and Wunsch 1978]。図1の蛋白質配列から3本の配列を選び、これに対し Dynamic Programming を適用し、計算機的に最適なアライメントを求める。3本の配列を6本の配列から選ぶ組み合わせは20通りあるので、20通りの計算機的に最適なアライメントが生成された。その内、生物学的に意味のあるアライメントは6通りのみであった。

```
17.6 : -----LLD-F-KL-HPCDQTKTGF-GCT-TY-FPFLQLQGQIINCSL-AKT-DRS-T-TMPXTT
M-MULV : -----LLD-FL-----RQ-LTE-SPLRS-LVLLERSSPPMMRQMLT-KRLLTOMAGD-YAA-SQS-----+---+---+
RSV : VASQGATGQ-FPLA-EAQQELT-ALGGQPLA-Q-HA-----CH13H0H-APVYQTCPC-BSA-FALAEQ-VI-
```

図2：計算機的に最適であるが生物学的には間違っているアライメントの例

図2に生物学的に間違っているアライメントの例を示す。この例では、RSV の配列で、2番目の共通な“H”を形成するべき“H”が、17.6 と M-MULV で、1番目の共通な“H”を形成するべき“H”と、同じカラムに並べられてしまっている。

### 3 Intelligent Refiner を用いたアライメントシステム

前の章で示したように、計算機的に最適なアライメントが、必ずしも、生物学的に意味のあるアライメントではない。以下、生物学的に意味のあるアライメントを作成するシステムを紹介する。我々のシステムは、Aligner と Intelligent Refiner の2つのモジュールにより構成される。

Aligner は、与えられた配列の計算機的に最適または準最適なアライメントを求める。これに Intelligent Refiner への入力となる。我々は、Aligner のアルゴリズムに[Ishikawa 1992]を採用しているが、他のアルゴリズムで作成されたアライメントを In-

telligent Refiner とすることができるので、以降、Aligner の説明は行わない。Intelligent Refiner は、Aligner により作成されたアライメントを生物学的な知識を用いながら検査に更新しながら、配列に含まれる重要な部位を特定していく。更新のサイクルが進むごとに、Intelligent Refiner は、配列が意味する生物学的情報をより詳細に把握していく。

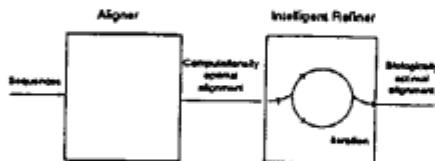


図3 : Intelligent Refiner を用いたアライメントシステム

図4にIntelligent Refinerの構成を示す。Refinement Rule Baseには、アライメントの更新を行う時に用いられるルールが登録されている。この中には、我々が生物学者から抽出したルールなどが含まれている。Control Moduleは、Refinement Rule Baseに登録されているルールを実行することによりアライメントを検査に更新していく。ルールを用いる時に必要な時にはBiological Knowledge Baseに蓄えられている生物学的知識を参照する。

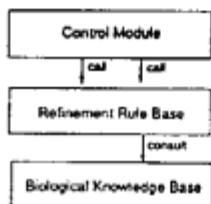


図4 : Intelligent Refiner の構成

Refinement Rule Baseには現在10箇のルールが登録されている。スペースの関係でこの内、2つのルールのみを図5に示す。また、Biological Knowledge Baseに登録されている知識の一部を図6に示す。モチーフの表現は[Bairoch 1991]に準拠している。

#### Rule 1

IF あるモチーフ ( $m_i$ ) がアライメントで特定され AND  
モチーフ  $m_i$  を持つ蛋白質が他のモチーフ  $m_j$  を持つ  
いれば

THEN Motif-finding routine が呼び出され、 $m_j$  を特定する。

#### Rule 2

IF Biological Knowledge Base に登録されているあるモチーフの中で、同一種類のアミノ酸  $x$  が ( $x_i$  and  $x_j$ ) 存在し、 AND

この2つのアミノ酸  $x_i$ 、  $x_j$  の間に他の保存アミノ酸が存在せず AND

refinement されるべきアライメントの中に、  $x$  が一部の配列を除き存在するカラム  $c_i$  と  $x$  が全ての配列に存在するカラム  $c_j$  があれば、

THEN Modification routine が呼び出され、以下の制約の下でアライメントを修正する（制約：  $c_i$  において  $x$  が存在しない配列を  $s_i$  とし、  $c_j$  の配列  $s_j$  に対応するアミノ酸を  $x_{j,i}$  とした時、  $x_{j,i}$  を  $c_i$  に編入する）。

図5 : Refinement Rule Base に登録されているルール

```

motif(name, zinc_finger, "H-I(3,5)-H-I(10,25)-C-X(3,5)-C").
motif(protein, kinase, "[LIV]-G-X-G-[FY]-[SG]-X-[LIV]").
motif(protein, kinase(tyrosine),
      "[LIVMFYC]-X-[HY]-X-D-[LIVMFY]-K-X(2)-N-[LIVMFY](3)".
upper_concept(kinase(tyrosine), kinase).
motif(protein, Protein, Motif) :-
    upper_concept(Protein, X), motif(protein, X, Motif).
  
```

図6 : Biologica Knowledge Base に登録されている知識

## 4 適用例

図2のIntelligent refinerを適用した(図7)。Rule 2などが適用され、生物学的に意味のある(Zinc Fingerを捕らえている)アライメントが作成された。なお、計算機的な評価値は計算機的に最適なアライメントは161であったが、生物学的に意味のあるアライメントでは156となり、評価は悪くなかった。これは、従来のアライメント手法では誤りとなる場合も、我々の手法を用いると適切なアライメントを作成できることを示している。

```

17.6 : SLDF--IQLVLLPQGKTHF--GET-TY-PPSOLLQIITIECTQIATVQVTPPTT-
R-HLYL : LLYV--LRL-LYDPSWQWLLRSQSPTRHQL-KLTETTKAGVVAASQVAVGTS-
RSV : YADQGQDQATPLAAGKLT ALKTQPAI--SKA---CH-ESQCD--RIVYVCPHQNSAPALEGTV-
(Evaluation value = 156)
  
```

図7 : Intelligent Refiner の適用例

## References

- [Bairoch 1991] Bairoch,A. Prosite : A dictionary of protein site and pattern : User manual Release 7.00, May 1991.
- [Berger and Manson 1991] Berger,M. and Manson,P. A novel randomized iterative strategy for aligning multiple protein sequences. Computer Application in the Biosciences, 7, 1991, pp.479-484.
- [Barton 1990] Barton,J.C. Protein Multiple Alignment and Flexible Pattern Matching. in Methods in Enzymology Vol.189, Academic Press, 626-645.
- [Hirosawa et al. 1992] Hirosawa,M., Ishikawa,M., Hoshida,M. Formulation of Protein Sequence Analysis using Knowledge 情報処理学会情報基礎研究会 ゲノム特集(1992)
- [Ishikawa et al. 1992] Ishikawa,M., Hoshida,M., Hirosawa,M., Toya,T. and Nitta,K. Protein Sequence Analysis by Parallel Inference Machine. Proc. Int. Conf. on Fifth Generation Computer Systems 1992.
- [Murata 1985] Murata,M. (1985) Simultaneous comparison of three protein sequences Proc. Natl. Acad. Sci. USA Vol. 82, 1985, pp.3073-3077.
- [Needleman and Wunsch 1970] Needleman,S.B. and Wunsch,C.D. (1970) A General Method Applicable to the Search for Similarities in the Amino Acid Sequences of Two Proteins. J. of Mol. Biol., 48, 443-453.