

近畿大学・農学部・生命情報学

ペアワイズアライメントと 配列相同性解析

2008年5月13日(火)

奈良先端大・情報・蛋白質機能予測学講座
川端 猛
takawaba@is.naist.jp

<http://isw3.naist.jp/IS/Kawabata-lab/home-ja.html>

授業予定			
日付	担当	講義	演習
4/8(火)	黒川	バイオインフォマティクス概論	
4/15(火)	黒川	配列解析1	IMCを使ったゲノム解析
4/22(火)	黒川	配列解析2	IMCを使った比較ゲノム解析
5/13(火)	川端	ペアワイズアライメントと配列相同性解析	
5/20(火)	川端	マルチプルアライメントと分子系統学基礎	配列相同性解析と系統樹作成演習
5/27(火)	川端	タンパク質配列の分類と機能推定	
6/3(火)	川端	タンパク質立体構造データの情報解析	タンパク質立体構造データの可視化演習
6/10(火)	川端	<試験>	
6/17(火)	金谷	ポストゲノム解析入門(トランスクリプトーム解析)	
6/24(火)	金谷	ポストゲノム解析入門(インタラクトーム解析)	発現プロファイル解析演習
7/1(火)	金谷	ポストゲノム解析入門(統合解析)	インタラクトーム解析演習・代謝物解析演習
7/8(火)	金谷	メタボローム解析(その1)	
7/15(火)	金谷	メタボローム解析(その2)	
7/22(火)	金谷	<試験>	

これから4回の講義の目標

イネのあるタンパク質のアミノ酸配列があったとして、

イネ: MAALSSAAVTIPSMAPSAPGRRRMRSSLV...

(1) 対応するほかの植物(たとえばマメ)のタンパク質を配列データベースから取り出したい

マメ: MATVTSTTBAIPSFSGGLKTNAATKVSAMA...

(2) どのアミノ酸とどのアミノ酸が対応するのか?

(3) もっとたくさんの似た配列があった場合、どれとどれが似ているのだろうか?

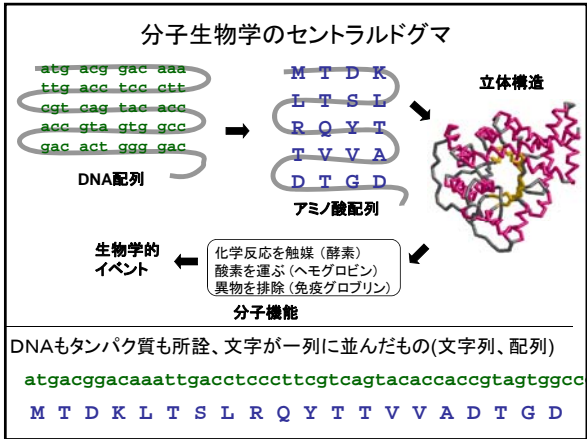
ポプラ: MAALSSAAVSVPSFAAATPMRSSRSRMV...

ナスナ: MAAITSATVTIPSFTGLKLAVSSKPKTLS...

(4) 機能的に大事なアミノ酸はどこだろう?

(5) どんな立体構造をしているのだろうか?

ペアワイズアライメント



「進化」とはDNAという文字列が変化すること

atgacggacaaaattgacctcccttcgctcagtagcacc

M T D K L T S L R Q Y T

atgacgaaacaaaattgacctcccttcgctcagtagcacc

M T N K L T S L R Q Y T

より正確には、個体のDNAが変化したあとに、その変異がその種の集団において定着する「集団遺伝学」的な過程が必要

- ① 個体のDNAに変異が生じる
- ② その変異が子孫に継承され、
- ③ 中立かつ正の淘汰が働けば、同じ変異を持った子孫が種の集団内で多数を占める

違う生物の同じ機能のタンパク質のアミノ酸配列

トリオースリン酸異性化酵素 (Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS))

>TPIS_HUMAN ヒト "Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)
 APSRKFFVGGNWKMNQRKQSLGELIGTLNAAKVPADTEVVCAPPT
 AYIDFARQKLDPKIAVAQAQNCYKVTNGAFTGEISPGMIKDCGATW
 VVLGHSERRHVFGESEDELIGQKVAHALAEGLGVIACIGEKLDERE
 AGITEKVVFEQTKVIADNVKDWKSVVLAIEPVAIGTGKTATPQQ
 AQEVHEKLRGWLKSNVSDAVAQSTRIIYGGSVTGATCKELASQPD
 VDGFVGGASLKPPEFVDIINAKQ

>TPIS_RABIT ウサギ "Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)
 APSRKFFVGGNWKMNQRKQSLGELIGTLNAAKVPADTEVVCAPPT
 AYIDFARQKLDPKIAVAQAQNCYKVTNGAFTGEISPGMIKDCGATW
 VVLGHSERRHVFGESEDELIGQKVAHALAEGLGVIACIGEKLDERE
 AGITEKVVFEQTKVIADNVKDWKSVVLAIEPVAIGTGKTATPQQ
 AQEVHEKLRGWLKSNVSDAVAQSTRIIYGGSVTGATCKELASQPD
 VDGFVGGASLKPPEFVDIINAKQ

違う生物の同じ機能のタンパク質のアミノ酸配列

トリオースリン酸異性化酵素 (Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS))

>TPIS_HUMAN ヒト "Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)
 APSRKFFVGGNWKMNQRKQSLGELIGTLNAAKVPADTEVVCAPPT
 AYIDFARQKLDPKIAVAQAQNCYKVTNGAFTGEISPGMIKDCGATW
 VVLGHSERRHVFGESEDELIGQKVAHALAEGLGVIACIGEKLDERE
 AGITEKVVFEQTKVIADNVKDWKSVVLAIEPVAIGTGKTATPQQ
 AQEVHEKLRGWLKSNVSDAVAQSTRIIYGGSVTGATCKELASQPD
 VDGFVGGASLKPPEFVDIINAKQ

>TPIS_YEAST 酵母 "Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)
 ARTFFVGGNFKLNGSKQSIKEIVERLNTASIPENVEVVICPPATY
 LDYSVSLVKKPQVTVGAQNAYLKASGAFTGENSVQIKDVGAKWV
 ILGHSERRSYFHEDDKFIADKTKFALGGVGVILCIGETLEKKA
 GKTLDVVERQLNAVLEEVKDWNTVNVVAYEPVVAIGTGKTATPQ
 QDIHASIRKFLASKLGDKAASELRILYGGSSANGSNAVTFKDKADV
 DGFVGGASLKPPEFVDIINSRN

違う生物の同じ機能のタンパク質のアミノ酸配列

トリオースリン酸異性化酵素 (Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS))

>TPIS_HUMAN ヒト "Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)
 APSRKFFVGGNWKMNQRKQSLGELIGTLNAAKVPADTEVVCAPPT
 AYIDFARQKLDPKIAVAQAQNCYKVTNGAFTGEISPGMIKDCGATW
 VVLGHSERRHVFGESEDELIGQKVAHALAEGLGVIACIGEKLDERE
 AGITEKVVFEQTKVIADNVKDWKSVVLAIEPVAIGTGKTATPQQ
 AQEVHEKLRGWLKSNVSDAVAQSTRIIYGGSVTGATCKELASQPD
 VDGFVGGASLKPPEFVDIINAKQ

>TPIS_ECOLI 大腸菌 "Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)
 MRHPLVMGNWKLNGSRHMVHELVSNLKELAGVAGCAVAIAPP
 YIDMAKREAEAGSHIMLGAQNVLDLNSGAFGTGETSAAMLKDIGA
 QYIIIGHSERRTYHKESDELIAKFAVLKEQGLTPVLCIGETEAE
 NAGKTEEVCARQIDAVLTKQGAFAFEGAVIAYEPVVAIGTGKSA
 TPAQAQAVHKFIRDHIAKVDANIAEQVVIQYGGSVNANAAELFA
 QPDIDGALVGGASLKDADAFVIVKAAEAAKQA

進化的なイベント: 置換 と 削除・挿入

トリオースリン酸異性化酵素 (Triosephosphate isomerase (EC 5.3.1.1) (TIM, TPIS)) の場合

ヒト (TPIS_HUMAN) とウサギ (TPIS_RABIT) の比較

```
HUMAN 1: APSRKFFVGGNWKMNQRKQSLGELIGTLNAAKVPADTEVVCAPPTAYIDFARQKLDPKIA:60
*****
RABIT 1: APSRKFFVGGNWKMNQRKQSLGELITLNAAKVPADTEVVCAPPTAYIDFARQKLDPKIA:60
*****
```

TPIS_HUMAN 248 vs TPIS_RABIT 248 SeqID 98.4 %

置換 (substitution) : アミノ酸・核酸の変化

ヒト (TPIS_HUMAN) と大腸菌 (TPIS_ECOLI) の比較

```
HUMAN 4: RRFVGGNWKMNQRKQSLGELIGTLNAAKVPADTEVVCAPPTAYIDFARQKLDPEKIAV:61
*****
ECOLI 2: RHPLVMGNWKLNGSRHMVHELVSNLKELAGVAGCAVAIAPPYIDMAKREAEAGSHIML:61
*****
```

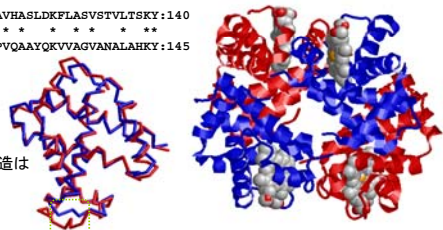
TPIS_HUMAN 248 vs TPIS_ECOLI 255 SeqID 45.9 %

削除・挿入 (insertion, deletion ; indel)

配列の類似と立体構造の類似

ヒトのヘモグロビンのα鎖とβ鎖 (SeqID 46.0%)

```
Alpha 2: LSPADKINVAANGKVGAGEYGAELERMFLSFPTTKTYFPHF-DLS-----HGSAQV:55
*****
Beta 3: LTPEEKSAVTALNGK--NVDVGGGALGRLLVVPWTRFFESFGDLSTPDVAVMNPVKV:60
*****
Alpha 56: KGHGKQVADALTNVAHVDDMPNALSALSDLHAHKLKRVDPVNFKLLSHCLLVTLAAHLPA:11
*****
Beta 61: KAHGKVLGAFSDGLAHLNLRKGTFAATLSLHCDKLRHVDPENFRLLGNVLCVLAHFGK:120
*****
Alpha 116: EFTPAVHASLDFKFLASVSTVLTLSKY:140
*****
Beta 121: EFTPFVQAAYQKVVAGVAMALAHKY:145
*****
```



機能や立体構造はよく似ている

配列の類似を知るとは立体構造予測につながる

配列比較 (配列相同性検索) の基本論理

① 2つの DNA / アミノ酸 の文字列が似ている



② 進化的に関係がある (相同) から似ている



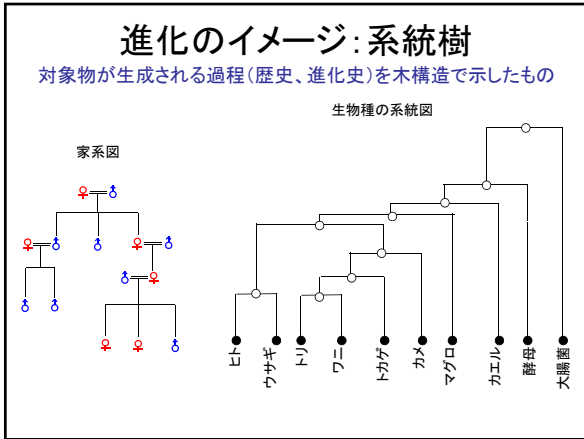
③ 進化的に関係があるなら、他の生物学的な性質 (機能、立体構造など) も似ているはず

相同性の発見により、他の生物学的な性質を予測できる

類似 (similarity)

相同 (homology): 進化的な原因によるもの。祖先を共有。
 (進化史の中である時点まで同じであったから似ている)

相似 (analogy): それ以外の原因によるもの



2つの配列を比較するには?

- 類似性のスコア関数の定義**
文字の間の類似性をどうやって定量するか?
 ACFDE
 ** * 3つ同じだから3点?
 ACEEE
 FとEの対応とDとEの対応は等価だろうか?
- アライメント**
どうやって文字と文字を対応づけるか?

ABCDEF	→	ABCDEF		BCDEF	→	-BCDEF-
CDE		--CDE--	***	ABEEFG		AB-EEFG

 もっと長いときはどうやって計算する?

スコア関数の定義

(1)一致・不一致スコア

$$S(A, B) = \begin{cases} \alpha & A = B \\ \beta & A \neq B \end{cases}$$

もっとも簡単。DNAの場合によく使われる。
BLASTの核酸のデフォルトは、 $\alpha=1, \beta=-3$

	A	T	G	C
A	1	-3	-3	-3
T	-3	1	-3	-3
G	-3	-3	1	-3
C	-3	-3	-3	1

問題点: 文字列間の類似性を捉えられない。
 L(ロイシン, 疎水性) → V(バリン, 疎水性) : 起こりやすい
 L(ロイシン, 疎水性) → E(グルタミン酸, 一荷電) : 起こりにくい

(2)対数オッズスコア(log odds score)

$$S(A, B) = \log \frac{P_{evo}(A, B)}{P_{rand}(A)P_{rand}(B)}$$

2つの異なるタンパク質のあるサイトのアミノ酸がA, Bであったとき、

Protein1: XXXXA
 Protein2: XXXXB

$P_{evo}(A, B)$: 進化的な関係からAとBの対応が生じた確率
 $P_{rand}(A) \cdot P_{rand}(B)$: 偶然にAとBの対応が生じた確率。

BLOSUM62 (blastpのデフォルトで使われている置換スコア行列)

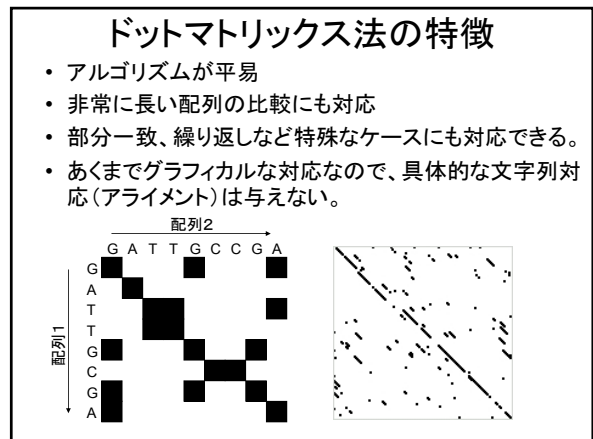
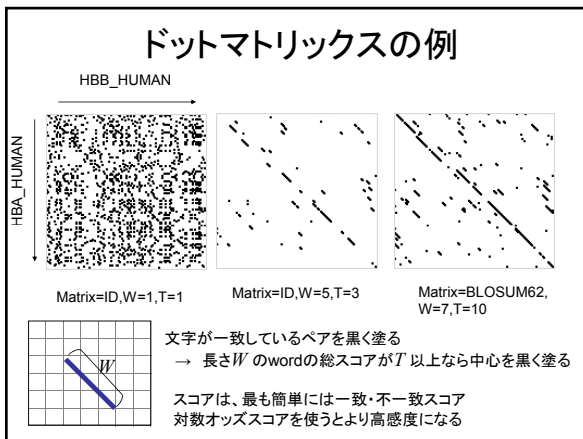
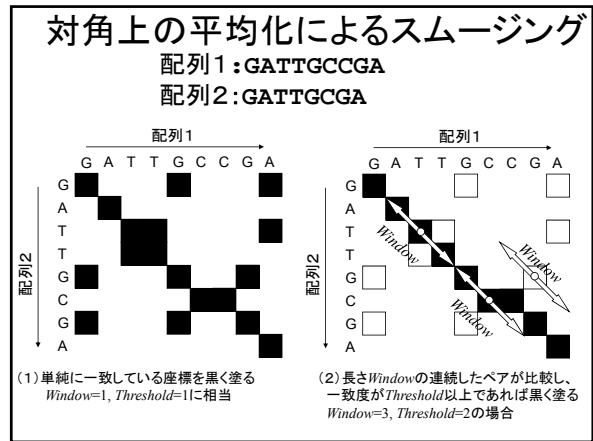
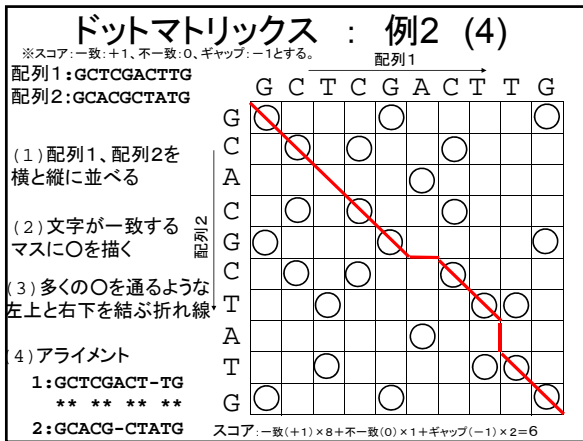
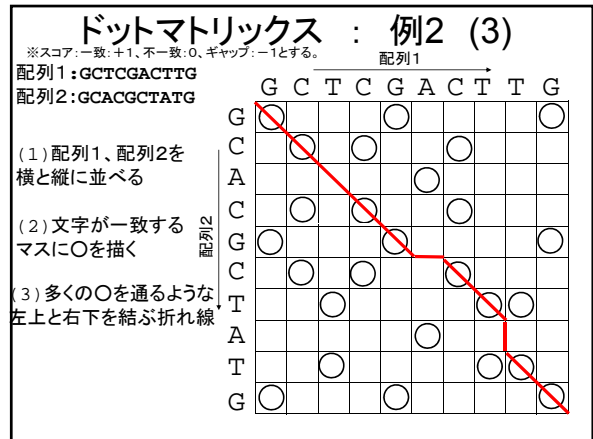
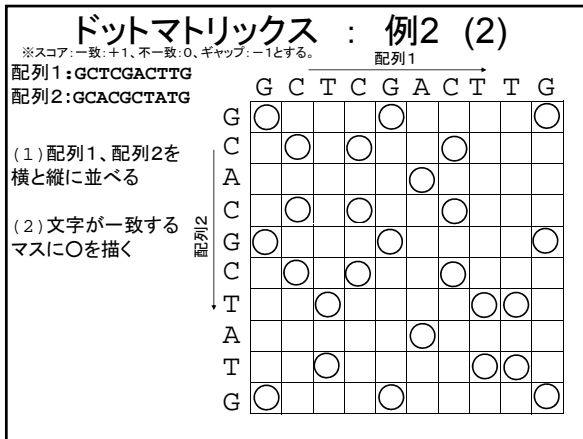
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	V	B	Z	X	*	
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-2	-1	1	0	-3	-2	0	-2	-1	0	-4	
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3	-1	0	-4	
N	-2	0	6	1	-3	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3	3	0	-1	-4	
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3	4	1	-4	
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1	-3	-3	-2	-4
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2	0	3	-1	-4
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2	1	4	-1	-4
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3	-1	-2	-1	-4
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	-2	-3	0	0	-1	-4
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	3	1	0	-3	-2	-1	-3	-1	3	-3	-3	-1	-4
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	2	2	0	-3	-2	-1	-2	-1	1	-4	-3	-1	-4
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	0	1	-1	-4	
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	1	5	0	-2	-1	-1	-1	-1	1	-3	-1	-1	-4
F	-2	-3	-3	-3	-2	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1	-3	-3	-1	-4	
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	4	7	-1	-1	-4	-3	-2	-2	-1	-2	-4
S	1	-1	1	0	-1	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2	0	0	0	-4	
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-2	-1	1	5	-2	-2	0	-1	-1	0	-4	
W	-3	-3	-4	-4	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3	-4	-3	-2	-4	
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	7	-1	-3	-2	-1	-4	
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4	-3	-2	-1	-4
B	-2	-1	3	4	-3	0	-1	0	-3	-4	0	-3	-3	-2	0	-1	-4	-3	-3	4	1	-1	-4	
Z	-1	0	0	1	-3	3	4	-2	0	-3	-3	1	-1	-3	-1	0	-1	-3	-2	1	4	-1	-4	
X	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-2	0	0	-2	-1	-1	-1	-1	-4	
*	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	1	

スコアの計算例

AFDC S(A,A) + S(F,E) + S(D,E) + S(C,C) = 12
 AECC 4 -3 2 9

ギャップがある場合はギャップのスコア(ギャップペナルティ)を設定する

AFDGC S(A,A) + S(F,E) + S(D,E) + gap + S(C,C) = 10
 AEE-C 4 -3 2 -2 9

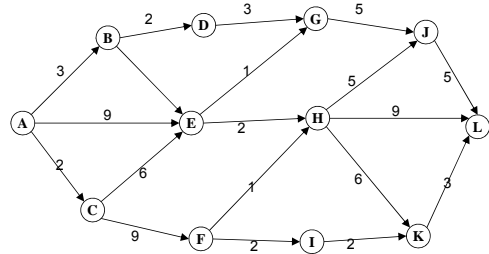


動的計画法によるアライメント

- アライメント問題は、有向グラフの最適経路問題と等価
- 有向グラフの最適経路問題は動的計画法 (Dynamic Programming) と呼ばれるアルゴリズムで解ける。
- $O(NM)$ の計算量 (文字列長の積に比例)

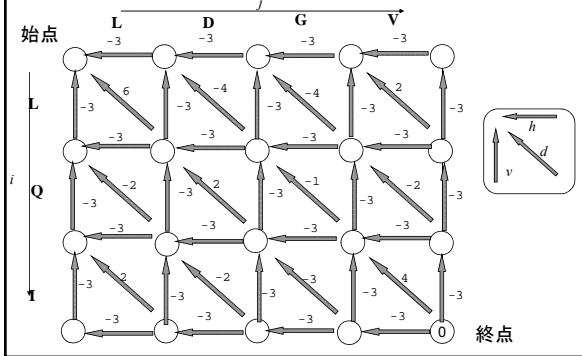
最適経路問題

始点Aから終点Lにいたるエッジの得点の合計が最大となる経路を探す



アライメントを最適経路問題として考える

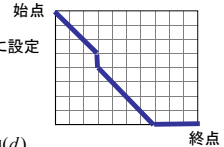
- 鉛直、水平に比較したい文字列を並べる
- 対角線のエッジには一致スコア、鉛直水平のエッジにはギャップスコアを書き込む
- 左上のノードから右下のノードへ至る最適経路を求める



グローバル・アライメントの解法 (Needleman & Wunsch, 1970)

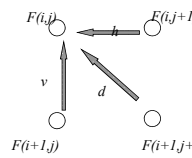
(0) 準備

右端の列、下端の行の格子点のスコアを0に設定



(1) 前向きステップ

$$F(i, j) = \max \begin{cases} F(i+1, j+1) + S(x_i, y_j) & \text{対角}(d) \\ F(i+1, j) + \text{Gap} & \text{鉛直}(v) \\ F(i, j+1) + \text{Gap} & \text{水平}(h) \end{cases}$$



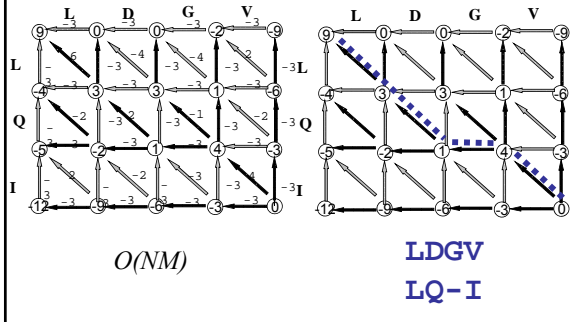
(2) 後ろ向きステップ

始点を起点にして進む。終点に到着したら終了。

動的計画法の手続き

(1) Forward

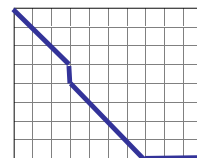
(2) TraceBack



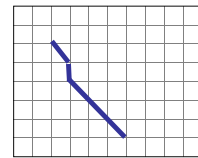
グローバルとローカルの格子上の違い

ACDEFGHKL M ACDEFGHK-LM FGHK-L
AFGHKKL A---FGHKKL- FGHKKL

グローバル ローカル



グローバル

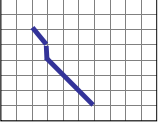


ローカル

ローカルアライメントの解法 (Smith & Waterman, 1981)

(0)準備
格子の端のスコアを0に設定

(1)前向きステップ

$$F(i, j) = \max \begin{cases} F(i+1, j+1) + s(x_i, y_j) & \text{対角}(d) \\ F(i+1, j) + \text{Gap} & \text{鉛直}(v) \\ F(i, j+1) + \text{Gap} & \text{水平}(h) \\ 0 & \text{終結}(0) \end{cases}$$



(2)後ろ向きステップ
最大のスコアのノードを探し、そのノードを起点にして辿る。パス'0'が現れたら終了

配列相同性検索

- BLASTを中心として -

配列相同性検索

→クエリ配列を配列データベースと比較、相同な配列を探す



- 機能未知遺伝子の機能予測(アノテーション)
機能既知の配列との類似→機能の類似を示唆
- 立体構造予測
構造既知の配列との類似→構造の類似を示唆
- 遺伝子発見
既知遺伝子と類似している領域の発見→遺伝子の存在を示唆

配列データベースの中からクエリ配列と類似したエントリを見つけるには？

→ 動的計画法を繰り返し実行すればよい

- いかに高速に計算を実行するか
動的計画法は $O(NM)$ の計算時間
1,000~100,000配列の検索には時間がかかる
→ 高度なヒューリスティック解法の導入
- どれだけ似ていれば意味があるのか？
何をもって類似性の指標とするのか
同一残基率(%), スコア?
→統計的有意性の判断の導入

BLASTのアライメントアルゴリズム

動的計画法を使わず、独自のヒューリスティックアルゴリズムを開発

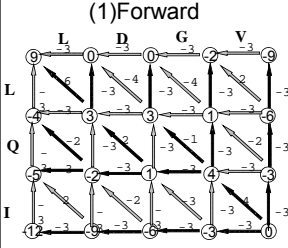
ヒューリスティック: 常に正しい解を返すわけではないが、多くの場合ままあの解を返すことが経験的に知られているアルゴリズム

153残基のクエリ配列を5977配列のデータベースと比較に要した時間(Pentium4)

私が書いたDP	16.989 sec
SSEARCH	2.911 sec
FASTA(ktup=1)	1.226 sec
FASTA(ktup=2)	0.608 sec
BLASTP	0.118 sec

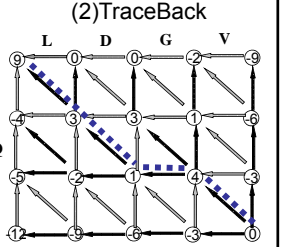
動的計画法の復習

(1)Forward



$O(NM)$

(2)TraceBack

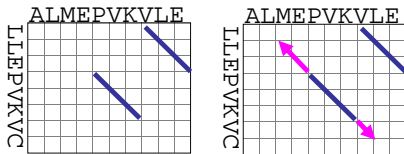


LDGV
LQ-I

BLASTのヒューリスティクス

目標: Smith&WatermanのローカルアライメントのDPの近似解

1. クエリの各wordに対し近隣wordのリストを作成
2. 近隣wordリストを用いてデータベースを検索
3. ヒットしたwordをgapで伸展(HSP)
4. さらにgap入りアライメントで伸展



BLASTP 2.2.1 [Apr-13-2001]

BLASTの出力例(1)

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

Query: RECA_ECOLI "Reca protein (Recombinase A)"
 (352 letters)
 Database: WgsdbP159nm
 3886 sequences: 705,110 total letters
 Searching.....done

Sequences producing significant alignments:	Score	E
	(bits)	Value
2reb-1 [c.37.1.11] RECA PROTEIN (E.C.3.4.99.37)	448	e-127
1g18A2 [d.48.1.1] RECA PROTEIN	70	9e-14
1g0uF [d.153.1.4] PROTEASOME COMPONENT C1	32	0.020
1byrA [d.136.1.1] ENDONUCLEASE	28	0.29
1g3a [c.37.1.10] CELL DIVISION INHIBITOR	28	0.38
1ct5A [c.1.6.2] YEAST HYPOTHETICAL PROTEIN, SELENOMET	28	0.49
1g0uB [d.153.1.4] PROTEASOME COMPONENT PUP2	27	1.1
1e32A2 [c.37.1.13] P97	26	1.4
1g0uA [d.153.1.4] PROTEASOME COMPONENT Y7	26	1.9
1cp2A [c.37.1.10] NITROGENASE IRON PROTEIN	26	1.9
1f3a [c.37.1.12] HYPOTHETICAL ABC TRANSPORTER ATP-BINDING PROTEIN	25	2.4
1g32B2 [d.133.1.1] CMB38 MONOKLIDE DEHYDROGENASE	25	3.2
1byyA [c.72.1.1] ADENOSINE KINASE	25	3.2
1skY83 [c.37.1.11] P1-ATPASE	25	3.2
1g6a [c.37.1.13] CAG-ALPHA	25	4.2
1cm5A [d.3.1.6] URQUHIN WH1-UBAL	24	7.1
184y- [c.93.1.1] L-ARABINOSE-BINDING PROTEIN (MUTANT WITH MET 1...	24	7.1
2tpa [c.1.3.1] THIAMIN PHOSPHATE SYNTHASE	24	7.1

BLASTの出力例(2)

```

>2reb-1 [c.37.1.11] RECA PROTEIN (E.C.3.4.99.37)
Length = 243
Score = 448 bits (1152), Expect = e-127
Identities = 243/266 (91%), Positives = 243/266 (91%), Gaps = 23/266 (8%)
Query: 3 DENKQKALAAALGQIEKQFKGSGIMRLGDRSMDVETISTGSLSDIALGAGLPMGRIV 62
Sbjct: 1 DENKQKALAAALGQIEKQFKGSGIMRLGDRSMDVETISTGSLSDIALGAGLPMGRIV 60
Query: 63 E1YGFESSGKTTLLQVIAAQRKGTCAFIDAEHALDPIYARKLGVDDINLCSQPDYG 122
E1YGFESSGKTTLLQVIAAQRKGTCAFIDAEHALDPIYARKLGVDDINLCSQPDYG
Sbjct: 61 E1YGFESSGKTTLLQVIAAQRKGTCAFIDAEHALDPIYARKLGVDDINLCSQPDYG 120
Query: 123 EQALEICDALARSQAVDVIIVDSVAALTPKAEIEGLAARMMSQMRKLAGNLL 182
EQALEICDALARSQAVDVIIVDSVAALTPKAEIEGLAARMMSQMRKLAGNLL
Sbjct: 121 EQALEICDALARSQAVDVIIVDSVAALTPKAEIEGLAARMMSQMRKLAGNLL 172
Query: 183 QQSNTLLIFINQIRMKIGVMPGNPTTGGNALKFYASVRLDIRRIGAVKSGENVVGSSET 242
QQSNTLLIFINQIRMKIGVMPGNPTTGGNALKFYASVRLDIRRIGAVKSGENVVGSSET
Sbjct: 173 QQSNTLLIFINQIRMKIGVMPGNPTTGGNALKFYASVRLDIRRIGAVKSGENVVGSSET 217
Query: 243 RVKVVKNKIAAPFRQARFQILYGEI 268
RVKVVKNKIAAPFRQARFQILYGEI
Sbjct: 218 RVKVVKNKIAAPFRQARFQILYGEI 243
>1g18A2 [d.48.1.1] RECA PROTEIN
Length = 60
Score = 70.1 bits (170), Expect = 9e-14
Identities = 30/56 (53%), Positives = 44/56 (78%)
Query: 272 GELVDLGVKLEKAGANYSYKGEIQGKANATWLNKDPETAKEIKKRELL 327
G L D + G V + L I K + G A N + Y + G E + + Q G K N A + L + N + A E I E K K + + E L
Sbjct: 192 K I I Y L A H E N K E D F E I S W C S L S E T N 219
    
```

BLASTの出力例(3)

```

Query: 243 RVKVVKNKIAAPFRQARFQILYGEI 268
Sbjct: 218 RVKVVKNKIAAPFRQARFQILYGEI 243
>1g18A2 [d.48.1.1] RECA PROTEIN
Length = 60
Score = 70.1 bits (170), Expect = 9e-14
Identities = 30/56 (53%), Positives = 44/56 (78%)
Query: 272 GELVDLGVKLEKAGANYSYKGEIQGKANATWLNKDPETAKEIKKRELL 327
G L D + G V + L I K + G A N + Y + G E + + Q G K N A + L + N + A E I E K K + + E L
Sbjct: 4 G S L I D M G V D Q G L I R K S G A N F Y E G E Q L G Q K E N A R N F L V E N A D V A D E I E K K I E K L 59
>1g0uF [d.153.1.4] PROTEASOME COMPONENT C1
Length = 242
Score = 32.3 bits (72), Expect = 0.020
Identities = 25/88 (28%), Positives = 47/88 (53%), Gaps = 9/88 (10%)
Query: 271 YGELVDLGVKLEKAGANYSYKGEIQGKANATWLNKDPETAKEIKKRELL 324
+G + G + +E +G++ YG G+G +A A L + +PE +A+E K+
Sbjct: 132 P G G V D K N G A H L Y L E P S G S Y W Y G K A A T G K G Q S A R A A E L K V D H P P G L S A R A V K Q A A 191
Query: 325 EL--LLSNPNSTPDFSVDDSE-GVAETN 349
++ L N D F ++ S ++ETN
Sbjct: 192 K I I Y L A H E N K E D F E I S W C S L S E T N 219
>1byrA [d.136.1.1] ENDONUCLEASE
Length = 152
Score = 28.5 bits (62), Expect = 0.29
    
```

どれだけ似ていれば意味があるのか？

類似性の指標

- 同一残基率(%)
直感的にわかりやすい。一般に30%ぐらいがしきい値とされる。感度が低く、アライメントの長さや不一致アアの類似性に鈍感

SLKA
* * 4/8 = 50 %
SELA Score = 4

SLKALLNKCKTFGWGAQ
* ** * * * * 8/16 = 50 %
SIRALDRRCKSFWAGKE Score = 55

- スコア
同一残基率より感度は高いが、比較する配列の長さに依存。長いほど高いスコアになる。

- E-value
スコアの統計的有意性。
ランダムな配列を比較した場合に、そのスコアが生じる可能性を見積もる。

E-value

E-value (expectation value)

ランダムな配列データベースを検索したときに、そのスコアS以上の値になるアライメントの本数の期待値

ランダムな配列とは: アミノ酸がランダムな順序に並んだ配列。ただし、アミノ酸の組成 → 平均的な値に従うとする
アミノ酸の長さ → 比較したアミノ酸の同じにする。

論理の流れ

ランダムな配列では起こりえないスコア
→ 偶然では起こりえないスコア → 進化的に関係がある類似性に違いない

値の大きさ

単位は本。小さいほどよく似ている。必ず0以上の値になる。

しきい値

原理的には1。経験的には0.0001から0.01ぐらい。

E-valueの計算に必要なパラメータ

$$E(S) = Kmn \cdot e^{-\lambda S}$$

- パラメータ定数 K, λ
→スコア行列とギャップペナルティに依存
 - m : クエリの残基長
データベースに含まれる全ての配列を一つにつなげた場合の長さ
 - n : データベースの残基長
- クエリ配列長とデータベースの大きさにE-valueは比例
比較した配列が同じでも、データベースのほかの配列の数が変わると、E-valueも変わってしまう。

```
BLASTP 2.2.1 [Apr-13-2001]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),
"Gapped BLAST and PSI-BLAST: a new generation of protein database search
programs", Nucleic Acids Res. 25:3389-3402.

Query= RECA_ECOLI "ReCA protein (Recombinase A)"
(352 letters)

Database: 40scopl.59nm
3886 sequences; 705,110 total letters

Searching.....done

Sequences producing significant alignments:

Score E
(bits) Value
2reb-1 [c.37.1.11] RECA PROTEIN (E.C.3.4.99.37) 448 e-127
lg18A2 [d.48.1.1] RECA PROTEIN 70 9e-14
lg0uF [d.153.1.4] PROTEASOME COMPONENT C1 32 0.020
lbyrA [d.136.1.1] ENDONUCLEASE 28 0.29
lg3qA [c.37.1.10] CELL DIVISION INHIBITOR 28 0.38
lct5A [c.1.6.2] YEAST HYPOTHETICAL PROTEIN, SELENOMET 28 0.49
lg0uD [d.153.1.4] PROTEASOME COMPONENT PUF2 27 1.1
le32A2 [c.37.1.13] P97 26 1.4
lg0uA [d.153.1.4] PROTEASOME COMPONENT V7 26 1.9
lg2A [c.37.1.10] NITROGENASE IRON PROTEIN 26 1.9
lf3oA [c.37.1.12] HYPOTHETICAL ABC TRANSPORTER ATP-BINDING PROTEIN 25 2.4
lgj2B2 [d.133.1.1] CARBON MONOXIDE DEHYDROGENASE 25 3.2
lgqvA [c.72.1.1] ADENOSINE KINASE 25 3.2
```

```
Query: 123 EQALEICDALARSGAVDIVVDSVAALTPKAEIEGEIGDSDHMLAARMMSQAMRKLGNL 182
EQALEICDALARSGAVDIVVDSVAALTPKAEIE GLAARMMSQAMRKLGNL
Sbjct: 121 EQALEICDALARSGAVDIVVDSVAALTPKAEIE-----GLAARMMSQAMRKLGNL 172

Query: 183 QKSNLLIFINQIRKMGVMPETTTGGNALKFYASVRLDIRRIGAVKEGENVVGSET 242
QKSNLLIFINQ TGGNALKFYASVRLDIRRIGAVKEGENVVGSET
Sbjct: 173 QKSNLLIFINQ-----TGGNALKFYASVRLDIRRIGAVKEGENVVGSET 217

Query: 243 RVKVVKNKIAAPFKQAEPQILYEGE 268
RVKVVKNKIAAPFKQAEPQILYEGE
Sbjct: 218 RVKVVKNKIAAPFKQAEPQILYEGE 243

Bit Score Raw Score
>lg18A2 [d.48.1.1] RECA PROTEIN
Length = 60
Score = 70.1 bits (1770), Expect = 9e-14
Identities = 30/56 (53%), Positives = 44/56 (78%)

Query: 272 GELVDLGVKELIEKAGAWYSYKGEKIGQKANATAWLKNDNPETAKEIKKVRRELL 327
G L+d+GV + LI K+GAW++Y+GE++GQGK NA +L +N + A EIEKK++E L
Sbjct: 4 GSLIDMGVDQGLIRKSGAWFTYEGEQLGQKKNARNFLVFNADVADEIEKKIEKEL 59

>lg0uF [d.153.1.4] PROTEASOME COMPONENT C1
Length = 242
Score = 32.3 bits (72), Expect = 0.020
Identities = 25/88 (28%), Positives = 47/88 (53%), Gaps = 9/88 (10%)

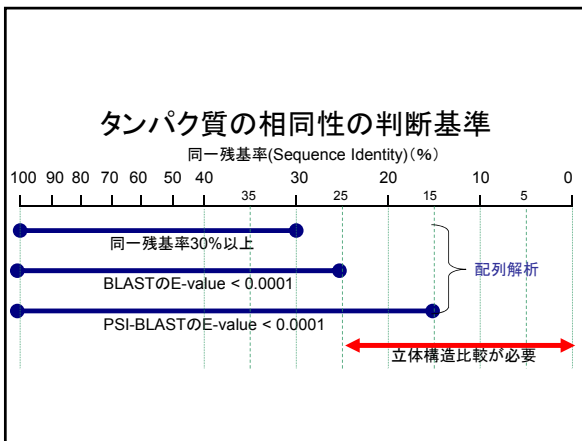
Query: 271 YGELVDLGVKELIEKAGAWYSYKGEKIGQKANATAWLK----DNPE--TAKEIKKVR 324
```

```
Database: 40scopl.59nm
Posted date: Jun 22, 2002 3:06 PM
Number of letters in database: 705,110
Number of sequences in database: 3886

Lambda K H
0.314 0.134 0.369

Gapped
Lambda K H
0.267 0.0410 0.140

Matrix: BLOSUM62
Gap Penalties: Existence: 11, Extension: 1
Number of Hits to DB: 469,543
Number of Sequences: 3886
Number of extensions: 18494
Number of successful extensions: 65
Number of sequences better than 10.0: 17
Number of HSP's better than 10.0 without gapping: 13
Number of HSP's successfully gapped in prelim test: 4
Number of HSP's that attempted gapping in prelim test: 50
Number of HSP's gapped (non-prelim): 17
length of query: 352
length of database: 705,110
effective HSP length: 79
effective length of query: 273
effective length of database: 398,116
effective search space: 108685668
effective search space used: 108685668
```



BLASTのプログラムの種類

	クエリ配列	データベース配列	比較回数	典型的な使用目的
blastn	核酸	核酸	2回 相補鎖にしたDB配列とも比較	ゲノムDNAのアンテセンス、cDNAのゲノムへのマッピング、非コーディング領域の比較
blastp	アミノ酸	アミノ酸	1回	タンパク質配列からの比較的遠縁のホモログの発見
blastx	核酸 (を翻訳したアミノ酸)	アミノ酸	6回 クエリから6通りのアミノ酸配列を生成して比較	ゲノムDNAから遺伝子 (タンパク質をコードしている領域) を発見する
tblastn	アミノ酸	核酸 (を翻訳したアミノ酸)	6回 クエリから6通りのアミノ酸配列を生成して比較	あるタンパク質をコードしているゲノムの領域を発見する
tblastn	核酸 (を翻訳したアミノ酸)	核酸 (を翻訳したアミノ酸)	36回 クエリ、DBとも6通りのアミノ酸配列を生成して比較	やや遠縁の生物種のゲノムを、その中にコードされたタンパク質で比較。DBに登録されていない遺伝子の発見を期待。

DNAには相補鎖があり、それぞれ3つのアミノ酸の読み枠がある

AGCTTTTTCATTCTGACTGCA
 |||
 TCGAAAAACAAGACTGACGT

DNAは二重らせん構造を作っているため、A⇔T、G⇔Cに入れ替えて、向きを逆にした相補鎖があるはず。

AGCTTTTTCATTCTGACTGCA
 S F S F X L Q
 A F H S D C
 L F I L T A

3つの核酸が1つのアミノ酸に翻訳されるので、読み枠をずらせば一本の核酸配列から3本のアミノ酸配列を作ることができる

※核酸よりアミノ酸で比較したほうがより遠縁のホモログを認識可能

blastp(アミノ酸対アミノ酸)によるタンパク質の機能予測

クエリ: *T.thermophilus*のタンパク質, データベース: 大腸菌の全タンパク質

```
BLASTP 2.2.3 [May-13-2002]
Query= X07_AAS80531.1 tche0 (144 letters)
Database: ecoli_ssa 4237 sequences; 1,350,094 total letters

Sequences producing significant alignments:
Score E
(bits) Value
infC NP_416233.1 "protein chain initiation factor IP-3" NC_000913 137 2e-34
rhaD NP_415030.1 "RhaD protein in RhaD element" NC_000913 28 0.19
pta NP_416800.1 "phosphotransacetylase" NC_000913 25 2.0
prsA NP_415725.1 "phosphoribosylpyrophosphate synthetase" NC_000913 25 2.7
yiaK NP_418032.1 "2,3-diketo-L-gulonate dehydrogenase, NADH-depe..." 24 3.5
ifh NP_417101.1 "4,5S-RNP protein, GTP-binding export factor, pa..." 24 4.6
yobB NP_415541.1 "putative dehydrogenase, NAD(P)-binding" NC_000913 24 4.6
ydfG NP_416057.1 "putative oxidoreductase" NC_000913 23 7.8

>infC NP_416233.1 "protein chain initiation factor IP-3" NC_000913
Length = 180

Score = 137 bits (346), Expect = 2e-34
Identities = 72/139 (51%), Positives = 92/139 (65%), Gaps = 1/139 (0%)

Query: 4 REALRLAQEMDLDLVVGPNADPPFARIMDYKWRVQOMKXXXXXXXXXXTEVKSIFR 63
REAL A+E -DLV + PRA+PPV RIMDI K+ YE+ +VK IKRF
Sbjct: 40 REALEKARANGYDLVEISFNAEPVCEIMDYKGRFLFKSKSKKQKQKQVIVQKIKRF 69

Query: 64 WKIDEDVQTKLGHKIFLQGHKRVVIMFGRREVAHPELGERILNRVTELDKLVAVVE 123
DE DYQ KE + RFL+DQ K K+* FGRRAAH ++G -LNRV +E+LAVVE
Sbjct: 100 PGTDRGDYQVKLSLIRFLREGDKARITLFRGRMAHQIGMEVLRNKKDLQELAVVE 159
```

blastp(アミノ酸対アミノ酸)の適用例)

ORFのアノテーション: *H.influenzae*のORF対大腸菌のORF

```
Query= HI0078_hinf0_AAC21753.1
Sequences producing significant alignments:
Score E
(bits) Value
cysS eco10 AAC73628.1 "cysteine tRNA synthetase" 730 0.0
metG eco10 AAC75175.1 "methionine tRNA synthetase" 39 5e-04
ileS eco10 AAC73137.1 "isoleucine tRNA synthetase" 39 0.001
leuS eco10 AAC73743.1 "leucine tRNA synthetase" 30 0.025
yidW eco10 AAC76718.1 "regulator protein for dgo operon" 28 1.3
→ HI0078はcysteine tRNA syntetase
```

```
Query= HI0083_hinf0_AAC21762.1
(71 letters)
Sequences producing significant alignments:
Score E
(bits) Value
ispB eco10 AAC76219.1 "octaprenyl diphosphate synthase" 23 3.1
lplA eco10 AAC77339.1 "lipoate-protein ligase A" 22 6.9
nlpA eco10 AAC76684.1 "lipoprotein-28" 22 6.9
bl372 eco10 AAC74454.1 "putative membrane protein" 22 6.9
mdaA eco10 AAC73938.1 "modulator of drug activity A" 22 9.0
→ HI0083は大腸菌にはホモログがない
```

参考文献

- 金久實 著 「ポストゲノム情報への招待」 (2001) 共立出版
- 中村保一他編 「バイオデータベースとウェブツールの手とり足とり活用法 改訂第2版」 (2007) 羊土社
- Arthur M.Lesk(岡崎康司、坊農秀雄 監訳)「バイオインフォマティクス基礎講義 一歩進んだ発想をみかのために」(2003), メディカル・サイエンス・インターナショナル
- D.W.Mount著、岡崎康司、坊農秀雄 監訳「バイオインフォマティクス-ゲノム配列から機能解析へ」第2版 メディカル・インターナショナル、2005年、11500円
- 阿久津達也 「バイオインフォマティクスの数理とアルゴリズム」(2007) 共立出版
- R.Durbin 他著、阿久津達也他訳 「バイオインフォマティクス - 確率モデルによる遺伝子解析」医学出版、2001年、9800円
- BLAST WEB page <http://www.ncbi.nlm.nih.gov/BLAST/>