

# IT IT-based Protein Sequence Analysis

Mi-Nyeong Hwang, JinSuk Kim  
Korea Institute of Science and Technology Information

< >  
가 가  
IT  
IT  
N-gram  
(Information Retrieval) (Text  
Categorization) 가  
:  
,

## < Abstract >

Methodology Though the sequence databases of proteins and DNAs are increasing in size exponentially, still exhaustive sequence search systems are commonly used in conducting biological researches. However, due to the advancement of IT information technologies, many information retrieval algorithms have been developed to search strings in large-scale text databases and are proved to be successful. We propose that these algorithms could also be applied to the biological data. N-gram indexing methods were applied to indexes from protein sequences of the protein sequences, and Information Retrieval algorithm and Text Categorization algorithm were applied to retrieve and to analyze for protein sequences.

*Keywords : Protein, Sequence Analysis, Information Retrieval*

## 1.

DNA  
DNA 가 (local alignment score)

가

가 가

(heuristic)

가

FASTA *k-tuple*

가 . FASTA *k-tuple*

2 *k* *k-tuple* *k-tuple*

[3].

2 가 *k-tuple*

가 BLAST(Basic Local Alignment Search Tool) 가

FASTA 1 n-gram

[4, 5]. CAFE , 2

(coarse search)

(inverted index) . (fine

search) 1 CAFE BLAST

," "

[6]. (full-text)

가

CAFE

가 DNA ,

[1].

n-gram ,

가 [2].

## 2. IT

### 2.1

C++ STL(Standard Template Libraries) [17]. Berkeley DB

, FASTA ProSeS(Protein Sequence Search)

B+ n-gram Berkeley DB

- 2.4GHz CPU, 3Gb 2.3

가 Ultra-160 SCSI

### 2.2

Williams PIR Super-family (global alignment) 가 (multiple alignment) PIR- [7, 8].

[9]. PIR-NREF

, ProSeS BLASTP [2002]. BLAST

BLASTP 가 [9].

. PIR-NREF

iProClass, INTERPRO, BLOCKS, PRINTS, PFAM, METAFAM, COG

SwissProt . PIR-NREF

2004 8 30 178 SwissProt 15

PIR-NREF ID PIR-ASDB, PIR-MIPS, PIRSUFAM, INTERPRO, PIRMOTIF, BLOCKS, PRINTS, PFAM, METAFAM, COG 가 . Swiss-prot

ProSLP(Protein Subcellular Localization Prediction)

가 1

```

@ProSeS
NREF_ID =NF00420179
REF_ID =GenPept:g414191; GenPept:g46348; GenPept:g39651209;
GenPept:g39648412; RefSeq:g39937351; RefSeq:g39934563;
PIR:S39073; SwissProt:LHB2_RHOPA;
TAXON_ID =1076
NAME =Light-harvesting protein B-800-850, beta chain B (Antenna pigment
protein, beta chain B) (LH II-B beta)
ORG =Rhodospseudomonas palustris CGA009
KEYWORDS =antenna complex; bacteriochlorophyll; inner membrane;
light-harvesting polypeptide; magnesium; transmembrane;
transmembrane protein;
PIR-ASDB =S39073;
PIR-MIPS =FAM0003052;
PIRSUFAM =SFO02900: Light-harvesting protein beta chain;
INTERPRO =IPRO00066: Antennacomplex, alpha/beta subunit;
IPRO02362: Antennacomplex, beta subunit;
PIRMOTIF =PCMO0969: PDC00748, Antenna complexes beta subunits signature
(PST: 17-48);
BLOCKS =IPB002362: Antenna complex, beta subunit;
PRINTS =PRO0674: LIGHTHARVSTB;
PFAM =PF00556: Antenna complex alpha/beta subunit (12-50);
METAFAM =m>S3907;
COG =
SLCC =Type II membrane protein. Inner membrane
SubCelLoc=720.01; 722.01; 722.02; 722.13;
ProSLP =720.01(100); 722.13(100);
LEN =51
SEQ =MADDPNKVWPTGLTIAESEELHKHVIDGTRIFGAIIVAHFLAYVYSPWLH

```

< 1 >

### 2.3 N-gram

, ( ) 가 .

file) [6, 7] DNA (inverted (string) ,  
 (heuristic) 가 , 가 n-gram . N-gram  
 가 OCR [11]. ,  
 가 n-gram  
 가 .  
 N-gram FASTA "k-tuple" BLAST "w-mer"  
 [13, 9]. n-gram , 가 n  
 ACEPITCH n 4 , n-gram ACEP, CEPI, EPIT,  
 PITC, ITCH가 .  
 1 n 3, 4, 5, 6 . ' . '

< 1 > 가 N-gram

N3-A20	3	20	8,000 (203)
N4-A20	4	20	160,000 (204)
N5-A20	5	20	3,200,000 (205)
N6-A18	6	18	34,012,224 (186)

PIR-NREF selenocysteine 가 21 .  
 X B, Z,  
 ..  
 (tri-gram)(N3-A20), (tetra-gram)(N4-A20), (penta-gram)(N5-A20) 20  
 (hexa-gram)(N6-18) 18  
 . 20 , 가 6 4 ,  
 . 20 2  
 . BLOSUM62 가 (V, I) (F,  
 Y) I Y . 18 3 4  
 가 , 가 [4].  
 N-gram Berkeley DB (searchable key)  
 (posting list)가 .  
 "ACEP" 36 (2), 127(3), 1074(1),  
 ACEP 36 , 127 , 1074  
 (posting  
 list) gzip .

## 2.4

가 , 가 가 . [7].  
 $q$   $d$  ( $Sim(q, d)$ )

$$Sim(q, d) = \frac{1}{W_d} \cdot \sum_{t \in q \wedge d} (w_{q,t} \cdot w_{d,t})$$

$$w_{q,t} = \log(f_{q,t} + 1) \cdot \log\left(\frac{N}{f_t} + 1\right)$$

$$w_{d,t} = \log(f_{d,t} + 1) \cdot \log\left(\frac{N}{f_t} + 1\right)$$

$$W_d = \log\left(1 + \sum_{t \in d} f_{d,t}\right)$$

$f_{s,t}$  s n-gram , N ,  $f_t$   
n-gram t가 ,  $w_{s,t}$  s t 가  
,  $W_d$  d .

## 2.5 가

가 PIR-NREF 100  
가 51 1609 316  
PIR-NREF BLAST  
BLAST , 1000  
, E-value 0.0001 . n-gram  
(reference test set) , ProSeS  
PIR-NREF 10,000 . N-gram  
가 BLAST (recall) (precision)  
[6]. 가가  
( $p$ ) BLAST

$$p = \frac{\text{BLAST sequences retrieved}}{\text{BLAST and non - BLAST sequences retrieved}}$$

( $r$ ) BLAST BLAST . ,

$$r = \frac{\text{BLAST sequences retrieved}}{\text{Total number of BLAST sequences}}$$

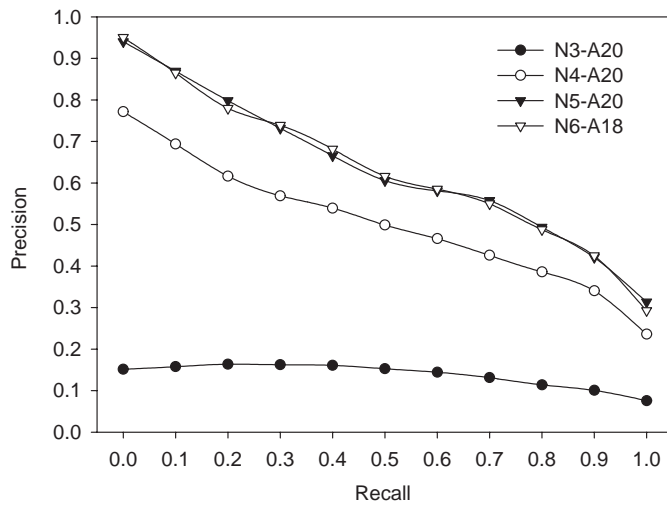
47가 n-gram 11 -  
0.0, 0.1, 0.2, , 1.0 11 1

### 3.

#### 3.1

1 47가 n-gram -  
(N6-A18) (N5-A20) 가 가  
(N3-A20) , (N4-A20) . N3-A20, N4-A20,  
N5-A20, N6-A18 n-gram 11 - 0.1376, 0.5038, 0.6342 0.6337  
11 - 0.25 0.40 , N5-A20  
N6-N18 0.63 [14].

가 BLAST



< 2> 가 n-gram 11 -

1 n-gram , (N6-A18) 가  
3가 18 가  
20 , 가  
가 1 2 - 가 2.8  
n-gram ,

가  
6

가 . n-gram

0.1 0.869 ProSeS 가

BLAST . 100 BLASTP

335 가 , 0.1 34 ,

ProSeS 30 ( 0.1 34 86.9%)

ProSeS BLASTP 0.1, 0.4, 0.7 1.0 가

73, 51, 38, 17 가 . ProSeS BLASTP가 0.1 73 , 0.4

51 가 . 9 0.1 0.1

가 97 , BLASTP

PIR-NREF 20 ProSeS 6

>NF01265541 Similar to 5'-nucleotidase, cytosolic III (Fragment) [Xenopus laevis]  
Length = 294

Score = 127 bits (320), Expect = 1e-28  
Identities = 96/288 (33%), Positives = 127/288 (44%), Gaps = 79/288 (27%)

Query: 15 PRALTDKIMTLIRDAGPSKFOVF-----PTP-----ISEQGDYAY 48  
P L DK+T I+ G K Q+ PT IS++G

Sbjct: 20 PEGLDQK ITR IQRGGQEKLQI ISDFDMLSRFSRNGERCPTCYNI IDNSNI ISDEGRK-- 77

Query: 49 DAKRQALYDHYHPLEISPVIPIDEKTKLMEIWIWKHELLIEGGLTYDAIKKSVANSSIA 108  
K + L+D Y+PLEI P I+EK LM EWW K H+L E + D + + V S

Sbjct: 78 --KLKCLFDIYYPLEIDPKKSIIEKYPLMVEIWSKAHDLFYEQRIQKDRLAQVVKESQAT 135

Query: 109 FREGVSELFFEFLEKKEIPVLIIFSAGLADVIEEVTLKSI SLELLLSYFCLYNEYAFVAYS 168  
R+G F L ++EIP+ IFSAG+ DV+EE

Sbjct: 136 LRDGYDLFFNSLYQREIPLFIIFSAGIGDVLEE----- 167

Query: 169 HSYQVLRQNLDRTFKNVKIVSNRMVFNDDGQLVSPKQKLIHVLNKNIEHALDMAAPLHDRL 228  
++RQ N K+VSN M F+D+G L FKG LIH NKN L

Sbjct: 168 ----IIRQ-AGVFHPNTKVVSNYMDFDNDGILTFGKGDLIHTYKNKSSVL----- 212

Query: 229 GVDIGEEDENVMKERRNVLMLMGDHLGDLRMSDGLD-YETRISIGFL 275  
+ E + R N+LL+GD LGDL M+DG+ E I IGF

Sbjct: 213 -----KDTEYFKEISHRTNILLGGDTLGDLTMDGVSTVENI IKIGFL 255

< 2> 97 BLAST "NF00667350 hypothetical protein  
At2g38680 [Arabidopsis thaliana]". ProSeS

2 97 BLASTP

ProSeS 2

129 IFSAG ProSeS

가

ProSeS,

3 97 BLASTP

ProSeS 129 IFSAG 256

GDLRM, ProSeS

2 2 3

ProSeS

2 3 BLAST , ProSeS 3  
 2 ProSeS BLAST  
 ProSeS

>NF01178199 10 days embryo whole body cDNA, RIKEN full-length enriched  
 library, clone:2610024B13 product:HSPC233 (PYRIMIDINE  
 5'-NUCLEOTIDASE) (EC 3.1.3.5) (URIDINE 5' MONOPHOSPHATE  
 HYDROLASE 1) (SIMILAR TO HYPOTHETICAL PROTEIN) homolog  
 [Mus musculus]  
 Length = 331

Score = 124 bits (310), Expect = 2e-27  
 Identities = 78/222 (35%), Positives = 118/222 (53%), Gaps = 49/222 (22%)

Query: 55 LYDHYHPLEISPVIPIDEKTKLMEEIHWGKTHELLIEGGLTYDAIKKSVANSSIAFREGVLS 114  
 L + Y+ +E+ PV+ ++EK M EW+ K+H LLIE G+ +K+ VA+S + +EG

Sbjct: 122 LKEQYYAIEVDPVLTVEEKFPYMWVEWYTKSHGLLIEQGIKAKLKEIVADSDVMLKEGYE 181  
 Query: 115 ELFEFLEKKEIPVLIIFSAGLADVIEEVTLSISLLELLSYFCCLYNEYAFVAYSHSYQVL 174  
 LF L++ IPV IFSAG+ DV+EEV ++ HS

Sbjct: 182 NLFCKLQQHGIPVFIIFSAGIGDVLEEVIRQA-----GVYHS---- 217  
 Query: 175 RQNLDRTFKNVKIVSNRMVFNDGQLVSFKGKLIHVLNKNHEALDMAAPLHDLRGVDIGE 234  
 NVK+VSN M F+++G L FKG+LIHV NK++ AL +

Sbjct: 218 -----NVKVVSNFMDFDENGVLKGFKGELIHVFNKHDGAL-----K 253  
 Query: 235 EDEENVNMKERRNVLIMGDHLGDLRMSDGL-DYETRISIGFL 275  
 + +K+ N++L+GD GDLM+DG+ + E + IG+L

< 3> 97 BLAST "NF01178199 10 days embryo  
 whole body cDNA [Mus musculus]". ProSeS

가

### 3.2

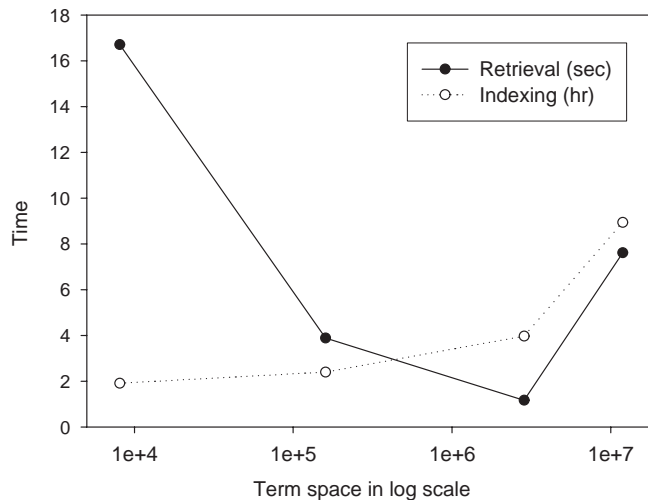
2 PIR-NREF , , 11  
 가 . PIR-NREF 1.41 release 404,532,594  
 가 3.7 14.0  
 (N5-A20) 7.6 가 , 가  
 ,  
 가  
 n-gram index stopping ProSeS  
 [15,16].

< > 가 BLAST , , , 11 - .



N-gram	(MB)	(hr)	( )	11
N3 - A20	1,514	1.91	16.70	0.1376
N4 - A20	3,125	2.40	3.88	0.5038
N5 - A20	3,075	3.97	1.17	0.6342
N6 - A18	5,655	8.94	7.61	0.6337
BLAST	583	0.05	44.10	-

4 n-gram n-gram  
가 ,  
C++ STL  
가



< 4> 가 n-gram (Term Space)  
N3-A20, N4-A20, N5-A20, N6-A18 8.0103, 1.6104, 2.8106, 1.2107

, N-gram 가 3 5 가 ,  
( 4). n 6 ,  
가 10 가  
1.17 가 ( 2). BLASTP 38  
가 가

4.

n-gram  
가 가 n-gram  
(gap)  
가 N-gram  
가 ProSeS, keyword suggestion  
subcellular localization, super-family classification 가  
(exhaustive comparison system)  
가 ProSeS가

[ ]

- [1] , , , "N- , Proceeding on 8  
, 3 , 171-182, 2004
- [2] , , "ProSLP: Penta-gram " , Proceeding on 8  
, 3 , 194-2004, 2004
- [3] Lipman, D. J. , Pearson. W. R. "Rapid and Sensitive Protein Similarity Searched" . *Science*, **227**, 1435-1441, 1985
- [4] Altschul, S. F. "Amino Acid Substitutions Matrices from an Information Theoretic Perspective" . *J. Mol. Biol.*, **219**, 555-665, 1991
- [5] Williams, H. E. and Zobel, J. "Indexing and Retrieval for Genomic Databases" . *IEEE Transactions on Knowledge and Data Engineering*, **14**, 63-78, 2002
- [6] Salton, G. and McGill, M. J. *Introduction to Modern Information Retrieval*, McGraw-Hill, New York, 1983
- [7] Witten, I. H., Moffat, A., and Bell, T. C. *Managing Gigabytes: Compressing and Indexing Documents and Images*, The Second Edition, Morgan Kaufmann Publishing, San Francisco, 1999
- [8] Wu, C. H., Huang, H., Yeh, L.-S. L., and Barker, W. C. Protein Family Classification and Functional Annotation. *Computational biology and Chemistry*, **27**, 37-4, 2003
- [9] Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Muller, W., and Lipman, D. J. Gapped BLAST and PSI-BLAST: A New Generation of Protein Database Search Programs. *Nucleic Acids Research*, **25**, 3389-3402, 1997
- [10] Wu, C., Huang, H., Arminski, L., CastroAlvear, J., Chen, Y., Hu, Z., Ledley, R. S., Lewis, K. C., Mewes, H. W., Orcutt, B. C., Suzek, B. E., Tsugita, A., Vinayaka, C. R., Yeh, L. S., Zhang, J. and Barker, W. C. The Protein Information Resource: An Integrated Public Resource of Functional Annotation of Proteins. *Nucleic Acids Research*, **30**, 35-37, 2002
- [11] Wilkinson, R. Chinese Document Retrieval at TREC-6. *NIST Special Publication 500-240: The Sixth Text Retrieval Conference (TREC 6)*, 25-29, 1998
- [12] Harding, S. M., Croft, W. B. and Weir, C. Probabilistic Retrieval of OCR Degraded Text Using N-grams. *European Conference on Digital Libraries*, 1997, 345-359, 1997

- [13] Lipman, D. J. and Pearson. W. R. Rapid and Sensitive Protein Similarity Searched. *Science*, **227**, 1435-1441, 1985
- [14] Oard, D. W and Gey, F. C. The TREC 2002 Arabic/English CLIR Track. *NIST Special Publication 500-251: The Eleventh Text REtrieval Conference (TREC 2002)*, 2002
- [15] Williams, H. E. and Zobel, J. Indexing Nucleotide Databases for Fast Query Evaluation. In *Peter M. G.Apers, Mokrane Bouzeghoub, and Georges Gardarin, Proc. International Conference on Advances in Database Technology (EDBT)*, Avignon, France, 275-288, 1996
- [16] Williams, H. E. and Zobel, J. Compression of Nucleotide Databases for Fast Searching. *Computer Applications in the Biosciences*, **13**, 549-554, 1997
- [17] Loverson, S. and Seltzer, M. Tree Houses and Real Houses: Research and Commercial Software, *Proceedings of the 2nd Workshop on Industrial Experiences with Systems Software (WIESS'02)*, Berkeley, CA: USENIX Association, December 2002, pp.55-6. 2002