

# N-gram Indexing for Protein Sequence Database

1, 1, 1, 1,2  
1

52

{jinsuk, mnhwang, snowy, ssh}@kisti.re.kr

2

28

가 가

4가 n-gram (tri-gram, tetra-gram, penta-gram, hexa-gram) PIR-NREF  
Penta-gram BLASTP 가  
38

## Abstract

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 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. Four n-gram indexing methods (tri-gram, tetra-gram, penta-gram, and hexa-gram) were applied to extract indexes from protein sequences of the PIR-NREF database, and their retrieval effectiveness and speed were measured. Penta-gram method showed the best results that its retrieval effectiveness matches for BLASTP and its retrieval speed was about 38 times faster than BLASTP program. Our protein sequence search service is accessible at <http://proses.kisti.re.kr>.<sup>1</sup>

---

<sup>1</sup>

1.

DNA

DNA

가

(local alignment score)

가

가

가

(heuristic)

가

FASTA *k-tuple*

*k-tuple*

가

FASTA

*k*

*k-tuple*

2

(Lipman *et al.*, 1985).

*k-tuple*

2

가

*k-tuple*

가 BLAST(Basic Local Alignment Search Tool)(Altschul, 1991)

가

FASTA 1

n-gram

(Williams *et al.*, 2002).

CAFE

2

(coarse search)

(inverted index)

(fine search)

1

CAFE

BLAST

(Salton *et al.*, Witten *et al.*, 1999).

(full-text)

가

CAFE

가

DNA  
 ,  
 n-gram  
 PIR-NREF  
 n-gram

## 2.

### 2.1

Berkeley DB  
 (Loverso *et al.*, 2002) C++ STL(Standard Template Libraries)  
 ProSeS(Protein Sequence Search) , FASTA  
 n-gram Berkeley DB  
 B+ N-gram 2.4  
 - 2.4GHz CPU, 3Gb Ultra-160 SCSI  
 가

### 2.2.

. Williams *et al.*(2002) PIR Super-family  
 . PIR- (global alignment)  
 가 (multiple alignment)  
 - (Wu *et al.*, 2003).  
 (Altschul *et al.*,  
 1997), PIR -  
 PIR-NREF (Wu *et al.*, 2002)  
 , ProSeS BLASTP . BLAST  
 BLASTP 가  
 (Altschul *et al.*, 1997). PIR-NREF  
 1.26 127  
 317 4 5

### 2.3 N-gram

, ( )

가 . , (inverted file) (Salton *et al.*, 1983, Witten *et al.*, 1999). DNA (string) , 가 (heuristic) 가 , 가 n-gram . N-gram 가 (Wilkinson, 1998) OCR (Harding *et al.*, 1997) , 가 , n-gram 가 . N-gram FASTA (Lipman *et al.*, 1985) “*k-tuple*” BLAST (Altschul *et al.*, 1997) “*w-mer*” . 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 (20 <sup>3</sup> )
N4 - A20	4	20	160,000 (20 <sup>4</sup> )
N5 - A20	5	20	3,200,000 (20 <sup>5</sup> )
N6 - A18	6	18	34,012,224 (18 <sup>6</sup> )

20 가 <sup>2</sup>. (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 (Altschul, 1991) 가

<sup>2</sup> 가 21 . PIR-NREF selenocysteine . B, Z, X .

(V, I) (F, Y) I Y 18  
 3 4 가 , 가 .  
 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 (Witten  
*et al.*, 1999).

$q$   $d$  ( $Sim(q, d)$ )

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

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

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

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

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

## 2.5

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)  
가  
(Salton *et al.*, 1983) . ( $p$ )  
BLAST . ,

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

( $r$ ) BLAST  
BLAST . ,

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

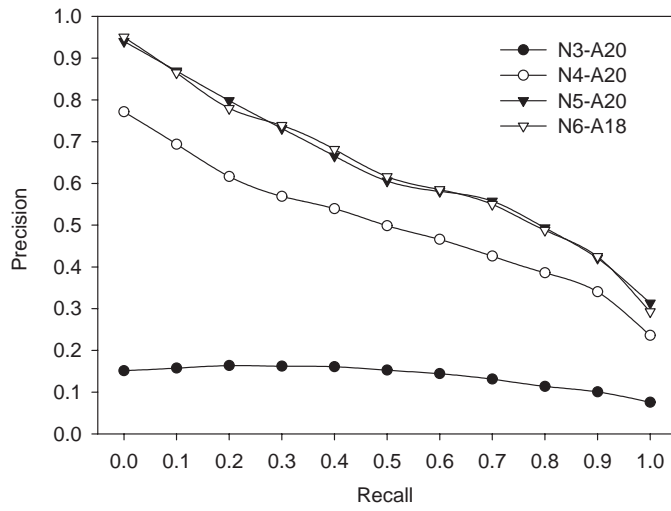
47가 n-gram 11 -  
. 0.0, 0.1, 0.2, ..., 1.0 11  
1 .

### 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  
(Voorhees, 2002, Oard *et al.*, 2002), N5-A20 N6-N18 0.63

가 BLAST .



1. 가 n-gram 11 -

1 가 n-gram (N6-A18) 3가 20

18 가 2 가 2.8

n-gram 6

0.1 0.869 ProSeS

가 BLAST 100

BLASTP 335 가 ProSeS

0.1 34 (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 PRALTDKMTLI RDAGPSKFQVF-----PTP-----ISEQGDAYY 48  
 P L DK+T I+ G K Q+ PT I S++G  
 Sbj ct: 20 PEGLODKI TRI QRGGQEKLOI I SDFDMTLSRFSRNGERCPTCYNI I DNSNI I SDEGRK-- 77

Query: 49 DAKRQALYDHYHPLEI SPVI PI DEKTKLMEEWVGKTHELLI EGGLTYDAI KKSIVANSSI A 108  
 K + L+D Y+PLEI P I+EK LM EWW K H+L E + D + + V S  
 Sbj ct: 78 --KLKCLFDI YYPLEI DPKKSI EEKYPMLVEWWSKAHDLFYEQRI QKDRLAQVVKESQAT 135

Query: 109 FREGVSELFEFLEKKEI PVLI FSAGLADVI EEVTLKSI SLLELLSYFCCLYNEYAFVAYS 168  
 R+G F L ++EIP+ IFSAG+ DV+EE  
 Sbj ct: 136 LRDGYDLFFNSLYQREI PLFI FSAGI GDVLEE----- 167

Query: 169 HSYQVLRQNLDRTFKNVKI VSNRMVFNDDGQLVSFKGKLI HVLNKNEHALDMAAPLHDRL 228  
 ++RQ N K+VSN M F+D+G L FKG LIH NKN L  
 Sbj ct: 168 ----I I RQ-AGVFHPNTKVVSNYMDFDNGI LTGFKGDLI HTYNKNSSVL----- 212

Query: 229 GVDI GEEDEENVNMKERRNVLLMGDHLGDLRMSDGLD-YETRI SI GFL 275  
 ++ E + R N+LL+GD LGDL M+DG+ E I IGFL  
 Sbj ct: 213 -----KDTEYFKEI SHRTNI LLLGDTLGDLTMADGVSTVENI I KI GFL 255

2. 97 BLAST "NF00667350

hypothetical protein At2g38680 [Arabidopsis thaliana]".

ProSeS

2 97 BLASTP

ProSeS

2 , 129 IFSAG

ProSeS

가

가

ProSeS,

, 3 97 BLASTP

ProSeS

129 IFSAG 256 GDLRM,

, ProSeS

2

2

3

ProSeS

2 3 BLASTP

, ProSeS

3

2

ProSeS BLASTP

, 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 LYDHYHPLEI SPVI PI DEKTKLMEEWGKTHELLI EGGLTYDAI KKS VANSSI AFREGVS 114  
 L + Y+ +E+ PV+ ++EK M EW+ K+H LLIE G+ +K+ VA+S + +EG  
 Sbj ct: 122 LKEQYYAI EVD PVL TVEE KFPYMVEWYTKSHGLLI EQGI PKAKLKEI VADSDV MLKEGYE 181

Query: 115 ELFEFLEKKEI PVL I FSAGLADVI EEVTLKSI SLLELLSYFCCLYNEYAFVAYSHSYQVL 174  
 LF L++ IPV IFSAG+ DV+EEV ++ HS  
 Sbj ct: 182 NLFGLKQQHGI PVFI FSAGI GDVLEEVI RQA-----GVYHS----- 217

Query: 175 RQNLDRTFKVKI VSNRMVFNDGQLV SFKGLI HVLNKNEHALDMAAPLHDLRGVDI GE 234  
 NVK+VSN M F+++G L FKG+LIHV NK++ AL +  
 Sbj ct: 218 -----NVK VSNFMDFDENGVLKGFKGELI HVF NKHDGAL-----K 253

Query: 235 EDEENVNMKERRNVLLMGDHLGDLRMSDGL-DYETRI SI GFL 275  
 + +K+ N++L+GD GD LRM+DG+ + E + IG+L  
 Sbj ct: 254 NTDYFSQLKDNSNI ILLGDSQGLRMDGVANVEHI LKI GYL 295

**3.97 BLAST . “NF01178199  
 10 days embryo whole body cDNA ... [Mus musculus]”. ProSeS**

가

**3.2**

2 PIR-NREF , ,

11 가 . PIR-NREF 1.26 release

404,532,594

가 3.7 14.0 . (N5-A20)

7.6 가 , 가 .

가

(Williams *et al.*, 1997)

n-gram *index stopping* (Williams *et al.*, 1996)

ProSeS

2. 가 BLAST , , , 11

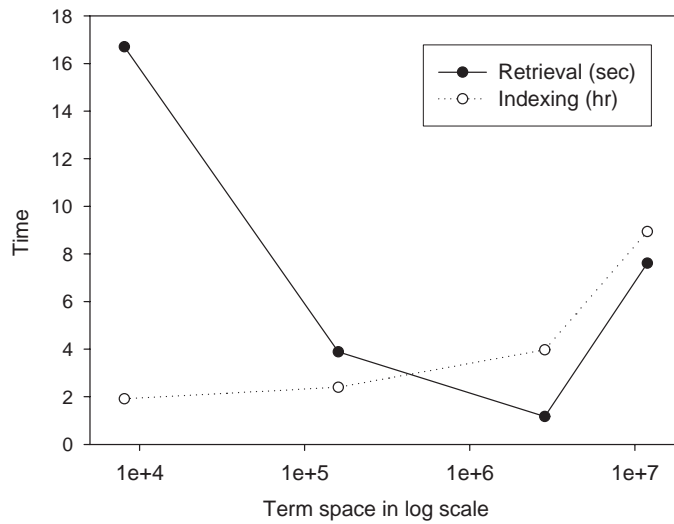
N-gram	(MB)	(hr)	( )	11
N3-A20	1,514	1.91	16.70	0.1376
N4-A20	3,125	2.40	3.88	0.5038
N5-A20	3,075	3.97	1.17	0.6342
N6-A18	5,655	8.94	7.61	0.6337
BLAST	583	0.05	44.10	-

4 n-gram

n-gram 가

C++ STL

가



4. 가 n-gram

(Term Space)

N3-A20, N4-A20, N5-A20, N6-A18

$8.0 \times 10^3$ ,  $1.6 \times 10^4$ ,  $2.8 \times 10^6$ ,  $1.2 \times 10^7$

가 , , . 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  
 가 . ProSeS, *keyword*  
*suggestion* *subcellular localization* , *super-*  
*family classification* 가  
 (exhaustive comparison system)  
 가  
 ProSeS가

Altschul, S. F. (1991) Amino Acid Substitutions Matrices from an Information Theoretic Perspective. *J. Mol. Biol.*, **219**, 555-665

Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Muller, W., and Lipman, D. J. (1997) Gapped BLAST and PSI-BLAST: A New Generation of Protein Database Search Programs. *Nucleic Acids Research*, **25**, 3389-3402

Harding, S. M., Croft, W. B. and Weir, C. (1997) Probabilistic Retrieval of OCR Degraded Text Using N-grams. *European Conference on Digital Libraries*, 1997, 345-359

Lipman, D. J. and Pearson. W. R. (1985) Rapid and Sensitive Protein Similarity Searched. *Science*, **227**, 1435-1441

- Loverson, S. and Seltzer, M. (2002) Tree Houses and Real Houses: Research and Commercial Software, *Proceedings of the 2<sup>nd</sup> Workshop on Industrial Experiences with Systems Software (WIESS'02)*, Berkeley, CA: USENIX Association, December 2002, pp.55-66
- Oard, D. W and Gey, F. C. (2002) The TREC 2002 Arabic/English CLIR Track. *NIST Special Publication 500-251: The Eleventh Text REtrieval Conference (TREC 2002)*, 2002
- Salton, G. and McGill, M. J. (1983) *Introduction to Modern Information Retrieval*, McGraw-Hill, New York, 1983
- Voorhees, E. M. (2002) Overview of TREC 2002. *NIST Special Publication 500-251: The Eleventh Text Retrieval Conference (TREC 2002)*, 2002
- Wilkinson, R. (1998) Chinese Document Retrieval at TREC-6. *NIST Special Publication 500-240: The Sixth Text Retrieval Conference (TREC 6)*, 25-29
- Williams, H. E. and Zobel, J. (1996) 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
- Williams, H. E. and Zobel, J. (1997) Compression of Nucleotide Databases for Fast Searching. *Computer Applications in the Biosciences*, **13**, 549-554
- Williams, H. E. and Zobel, J. (2002) Indexing and Retrieval for Genomic Databases. *IEEE Transactions on Knowledge and Data Engineering*, **14**, 63-78
- Witten, I. H., Moffat, A., and Bell, T. C. (1999) *Managing Gigabytes: Compressing and Indexing Documents and Images*, The Second Edition, Morgan Kaufmann Publishing, San Francisco, 1999
- 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. (2002) The Protein Information Resource: An Integrated Public Resource of Functional Annotation of Proteins. *Nucleic Acids Research*, **30**, 35-37
- Wu, C. H., Huang, H., Yeh, L.-S. L., and Barker, W. C. (2003) Protein Family Classification and Functional Annotation. *Computational biology and Chemistry*, **27**, 37-47