

DNA

An Efficient Homology Search Method for DNA Sequence Databases

(Sang-Kyoon Hong)*, (Jung-Im Won)**, (Jee-Hee Yoon)*,
 (Sang-Hyun Park)**, (Sang-Wook Kim)***

first) DNA (breadth-
 가 가
 : DNA

ABSTRACT

In molecular biology, DNA sequence searching is one of the most crucial operations. In this paper, we suggest an efficient homology search method for DNA sequence databases. The proposed index consists of two parts: the primary part represents the trie as bit strings without any pointers, and the secondary part helps fast accesses of the leaf nodes of the trie that need to be accessed for post processing. We also suggest an efficient approximate subsequence match algorithm based on that index. Our method employs a dynamic programming search technique which is driven by traversing the trie index in the breadth-first order. To verify the superiority of the proposed approach, we conducted a performance evaluation via a series of experiments. The results revealed that the proposed approach, which requires smaller storage space, can be a few orders of magnitude faster than the suffix tree.

Key words : DNA Sequence Database, Indexing, Trie

(R04-2003-000-10048-0)

IT

*
 **

+ : 2005 2 24 , : 2005 4 14

1.

DNA

가 , DNA

DNA

A, C, G, T 가

(suffix tree)[23] DNA

DNA

가

DNA

[8].

DNA

DNA

(homology search)

[24, 27]. DNA

S

[6, 11, 15, 25].

Q,

(tolerance) T가

, Q

Q,

가 T

S, S

T 가

Q, S,

가

(similarity function)

[18]

(edit distance)

T

DNA

[22]. BLAST

[1, 2] DNA

가

. BLAST

[9, 14, 18].

Smith-Waterman

[11] , 286M DNA

[22]

가

19G

:

가

[20].

가

DBMS

:

가

가

DBMS

(seamless integration)

2.

[23, 26].

DNA

DNA

Smith-

Waterman(S-W)

[22]

S-W

(trie)[10, 23]

(dynamic programming)

S Q

가

, DNA

(optimal local alignment)

가

($O(|Q \times S|)$)

가

BLAST[1, 2]

DNA

가

가 (near

optimal)

BLAST

가

가

가

DNA

DNA

(suffix tree)[23]

DNA

(inverted index)

[5, 7, 27],

[12],

[11, 17, 20]

[11]

(partitioned suffix tree)

(inverted index)

[20]

[5, 7, 27]

[16]

[17]

(interval)

가

A*-

[13]

(

)

[27]

가

가

BLAST

가

가

[12]

(wavelet)

가

DBMS

가

3.

DNA

DNA

(edit distance)

, global alignment

| | |
|----|-----|
| \$ | 000 |
| A | 001 |
| C | 010 |
| G | 011 |
| N | 100 |
| T | 101 |
| S | 110 |
| Y | 111 |

< 1>

| | |
|------------|-----------------|
| S1: ACGT\$ | 001010011101000 |
| CGT\$ | 010011101000 |
| GT\$ | 011101000 |
| T\$ | 101000 |
| S2: ACT\$ | 001010101000 |
| CT\$ | 010101000 |
| T\$ | 101000 |

< 2> DNA

DNA A, C, G, T
 가 11
 가 N
 A, C, G, T
 , B A가 C, G, T
 DNA 7
 가 , 3
 DNA
 가
 ,
 < 1> DNA
 '\$'
 , S1='ACGT' S2='ACT'
 DNA

<

1> 3
 < 2>
 DNA
 Algorithm 1
 , Algorithm 1 DNA
 가 .
 . 1
 가
 가
 가 ,
 가 . 2
 . < 3> < 2>
 S1 'ACGT\$'
 (001010011101000)

,

Algorithm 1: Index construction

Input : set of binary suffix sequences S

Output : binary suffix trie I

- 1 : S
- 2 :

2

가

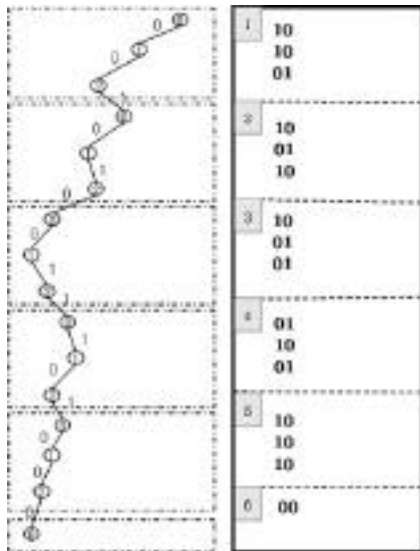
가 가

< 4> S2

'ACT\$(001010101000)' 가

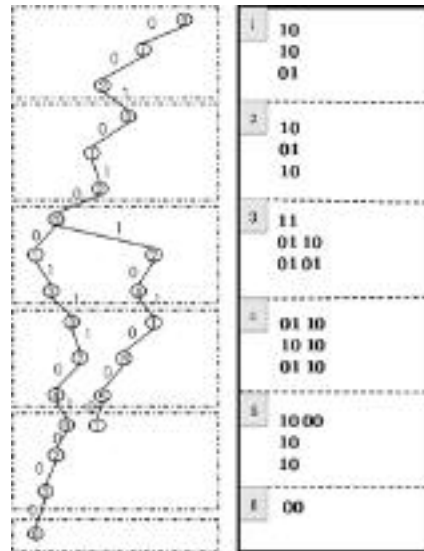
가 ,

2



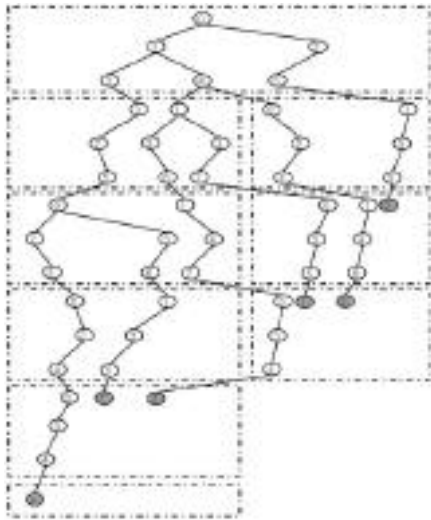
< 3>

S1 = 'ACGT\$'



< 4>

S2 = 'ACT\$'



< 5>

. < 5> < 2>

, < 6> < 5>

Algorithm 1

가
. < 3, 4, 5, 6>

. < 3, 4, 5>

| | | |
|---|-------------------------------|------------------------------|
| 1 | 11 11 10 01 11 01 | |
| 2 | 10 11 01 01 10 10 01 01 | 3 01 10 10 10 01 10 |
| 4 | 11 01 01 10 10 01 01 01 | 5 10 10 00 10 10 10 10 |
| 6 | 01 10 10 10 01 10 | 7 10 00 00 10 10 |
| 8 | 10 00 00 10 10 | |
| 9 | 00 | |

< 6>

16 가 ,

3 , 8 가 ,

. < 3> 2

3

가 가

< 5>

가

가

| # Page | Top | Bottom | Node | Addr |
|--------|-----|--------|------|------|
| 1 | 0 | 0 | 6 | 84 |
| 2 | 0 | 0 | 8 | 54 |
| 3 | 2 | 3 | 6 | 108 |
| 4 | 0 | 0 | 8 | 24 |
| 5 | 2 | 3 | 7 | 132 |
| 6 | 0 | 0 | 6 | 0 |
| 7 | 2 | 2 | 5 | 159 |
| 8 | 0 | 0 | 5 | 180 |
| 9 | 0 | 0 | 1 | 201 |

< 7>

3.2

가 . ,) (

가

가 가 2

4

< 6>

< 7>

가 # page

, Top

, Bottom

, Node

, Addr

가 ,

가 ,

DP $|Q|+1$ $S \ Q$
 2 DP $|S|+1$
 (optimal distance)

가
 (edit distance)
 [11, 23].

가

k

k

k

DP

0 k

k

[11, 17, 18, 25].

k

k

k

2

k

4.

3

3

4.1

'AGG'

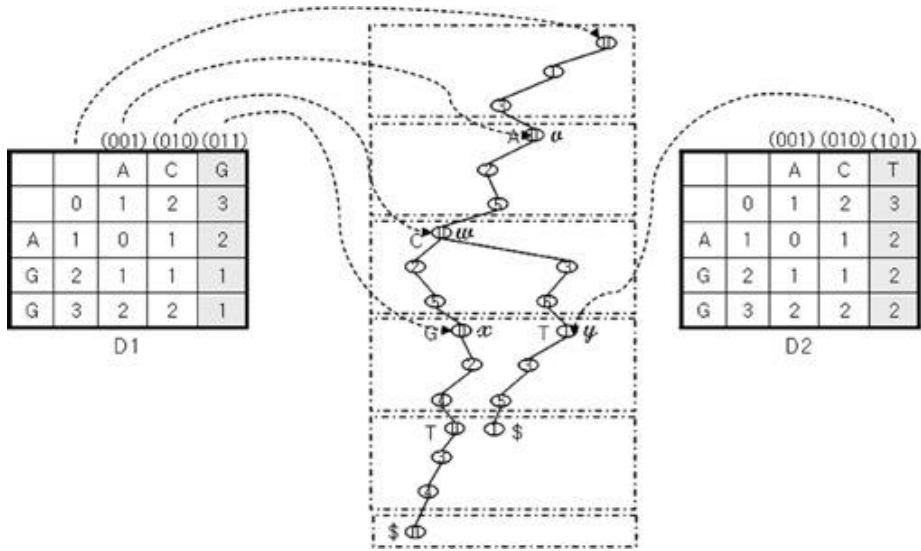
k가1

가

(DP)

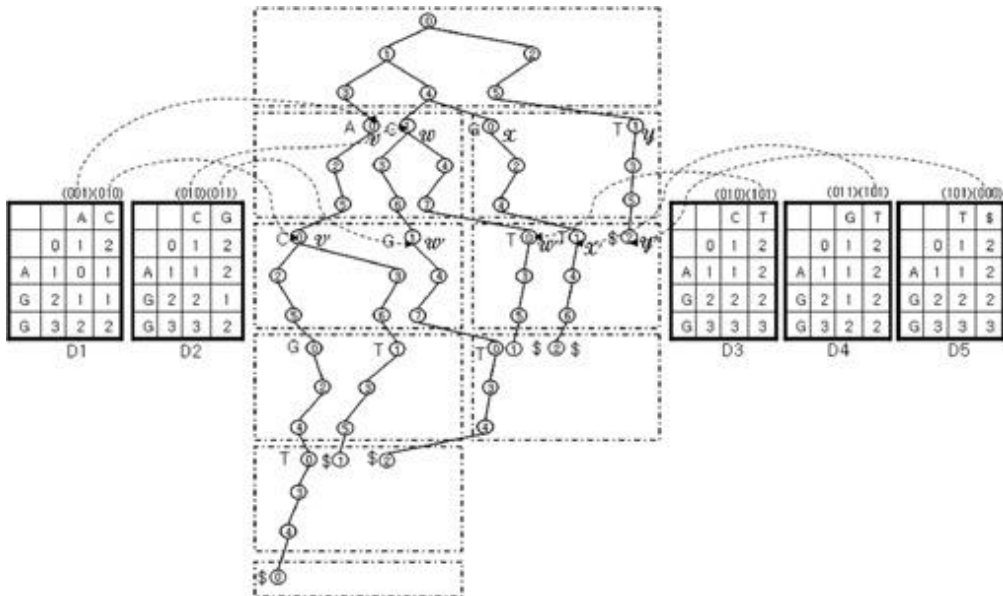
8

3



< 8 >

가 y . D2 y
 3 'T(101)'
 DP 가 DP
 , < 8 > v, w, x D2 'A','C'
 'A(001)'; 'C(010)'; D1 'A','C'
 'G(011)' 가 DP 가 D2 'T'
 D1 . DP 1
 가 가1 가
 가 ,
 가 'CGT(10011101)'
 가 . <
 8> D1 'G' 3.1
 1 가1 2
 가 x
 가
 , w



< 9>

‘CT(010101)’; ‘GT(011101)’; ‘T\$(101000)’
 5 DP D1, D2, D3, D4,
 D5 D3
 ‘C’ D2 ‘C’

가

< 8>

< 9>

가

v, w, x,

Search-Trie

y

4 DP

v, w, w, x, y,

l

Q

‘AC(001010)’; ‘CG(010011)’;

Q

, Q_node
 Q_pagenumber ,
 Q_node
 Enqueue (Line 1-2).
 while
 while (Line 4-33) Line 4
 Q_pagenumber
 Dequeue while
 (Line 7-21) Line 7-8 Dequeue
 Q_node
 Dequeue
 current_Node
 current_Node
 (Line 11-21). Line 11
 moreVisit
 TRUE . Line 12
 AppendBitString() CNI
 가
 current_Node
 current_Node
 가 .
 가 3 (Line14),
 가 가
 DP N 가 ,
 가 . Line 15
 AddColumn(currentDPT, newAdded_path)
 , current_Node
 DP current_DPT DP
 가
 DP

dist (Line 16).
 dist k
 가
 FindAnswer()
 CNI
 (
 moreVisit FALSE
 dist k
 FurtherVisit()
 moreVisit
 Q_node Enqueue
 TRUE
 Line 23-27
 가
 Q_pagenumber Enqueue ,
 Q_node
 Enqueue .
 4.2
 4.1 Search-Trie
 ,
 가 ,
 N
 N
 가
 Find_Answer() 가
 . 3.2

가 .

GenBank[19]

human chromosome 18, 19, 21

N ,
3가 ,

DNA

= {A, C, G, T}

가 , 가 N, S, Y
가 .

DNA 7

가 .

p 가 k
: p 0

' \$ ' 8가

k- p0 , p 1
k- p1

Windows 2003 Server

, 1GB

p1

200GB

Pentium IV 3.2GHz

PC

p 가 k

가

:

p k- 가

. Trie

p 가 k
: p k-

pt

pt 가

R*- [3]

R*- Maryland

Faloutsos

R*-tree

Version 2.0

2

5. 가

Suffix

. Suffix

[4]

incremental disk-based

Suffix 24%

32

1 () 2 ()

1 Trie

Suffix 가 . < 2 Trie

10>

4K 가 k

Suffix

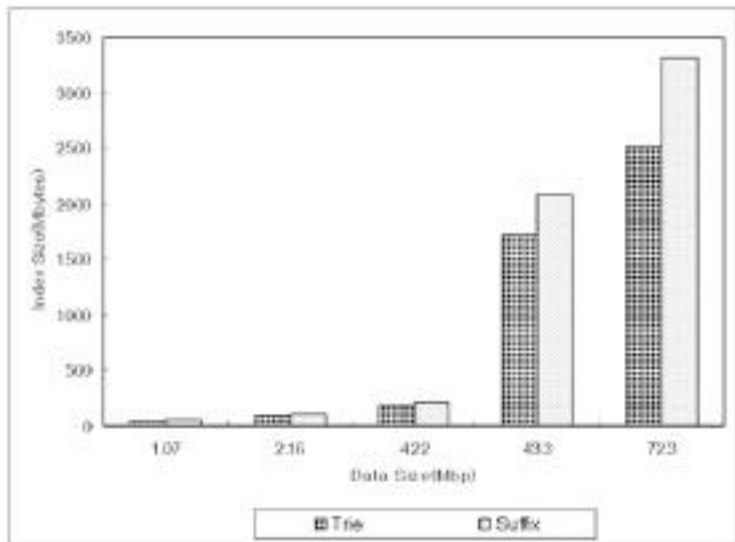
. Trie 43.3Mbp human chromosome 21

R- Trie

Suffix, Trie Trie-Depth

가 가 Trie-Breadth

Trie < 1> 가1



< 1> 1

| Query Length | # Result | Query Processing Time(msec) | | |
|--------------|----------|-----------------------------|----------------|---------------|
| | | Suffix | Trie-Depth | Trie-Breadth |
| 6 | 388,321 | 9962.2 (6074.2) | 1152.2 (488.7) | 980.2 (470.2) |
| 8 | 33,422 | 1800.5 (517.2) | 704.1 (278.5) | 333.5 (272.3) |
| 10 | 3,966 | 1301.4 (61.2) | 931.7 (217) | 264 (212.7) |
| 15 | 22 | 1224.4 (0) | 1062.8 (3.6) | 74.7 (2.3) |
| 30 | 4.1 | 1174.2 (0.1) | 1096.8 (2.6) | 90.7 (1.9) |

< 2> 2

| Query Length | # Result | Query Processing Time(msec) | | |
|--------------|-----------|-----------------------------|-----------------|----------------|
| | | Suffix | Trie-Depth | Trie-Breadth |
| 6 | 3,772,995 | 175668.2(78113.3) | 11048.8(4835.5) | 9640.9 (4778) |
| 8 | 524,469 | 22829.8 (8855.5) | 10863.4 (5140) | 6040.8 (5178) |
| 10 | 80,740 | 4043.9 (1185) | 16076.7 (4412) | 4891.2(4313.5) |
| 15 | 296 | 2808.5 (1.9) | 19332.2 (54.5) | 889.6 (49.2) |
| 30 | 18 | 2816.3 (0) | 32299.5 (5) | 938.6 (5.2) |

가 2 , < 2> 가
 , Trie R-
 ,
 Trie-Depth ,
 가 가
 가
 [18] 가
 Trie-Breadth 가
 가
 Suffix 가

| | | | |
|--------------|---------|----|---|
| | 가 | 가 | . |
| Trie-Breadth | , | | 가 |
| | Suffix, | | 3 |
| Trie-Depth | | 18 | |

6.

DNA

[1] S. Altschul, W. Gish, W. Miller, E. Myers, and D. Lipman, "Basic local alignment search tool," *Journal of Molecular Biology*, Vol. 215, No. 3, pp. 403-410, 1990.

[2] S. Altschul, T. Madden, A. Schaffer, J. Zhang, W. Miller, and D. Lipman, "Gapped BLAST and PSI-BLAST: A New Generation of Protein Database Search Programs," *Nucleic Acids Research*, Vol 25, No. 17, pp. 3389-3402, 1997.

[3] N. Beckmann, H. Kriegel, R. Schneider, and B. Seeger, "The R*-tree: An efficient and robust access method for points and rectangles," In *Proceedings of ACM SIGMOD International Conference on Management of Data*, pp. 322-331, 1990.

[4] P. Bieganski, J. Riedl, J. V. Carlis, "Generalized suffix trees for biological sequence data: applications and implementation," In *Proceedings of Hawaii International Conference on System Sciences*, Vol. 5, pp. 35-44, 1994.

[5] A. Califano and I. Rigoutso, "FLASH:

가 가

DBMS

(breadth-first)

- A Fast Look-up Algorithm for String Homology," In Proceedings of Intelligent System Conference for Molecular Biology, pp. 56-64, 1993.
- [6] A. L. Delcher, S. Kasif, R. D. Fleischmann, and J. Peterson, O. White, and S. L. Salzberg, "Alignment of whole genomes," *Nucleic Acids Research*, 27, pp. 2369-2376, 1999.
- [7] C. Fondrat and P. Dessen, "A Rapid Access Motif database (RAMdb) with a search algorithm for the retrieval patterns in nucleic acids or proteun databanks," *Computer Applications in the Biosciences*. Vol. 11, No.3, pp. 273-279, 1995.
- [8] C. Gibas and P. Jambeck, *Developing Bioinformatics Computer Skills*, O Reilly and Associates Inc., 2001.
- [9] R. Giegerich, S. Kurtz, and J. Stoye, "Efficient Implementation of Lazy Suffix Trees," *Softw. Pract. Exp.*, Vol 33, pp. 1035-1049, 2003.
- [10] E. Horowitz, S. Sahni, and S. Anderson-Freed, *Fundamentals of Data Structures in C*, Computer Science Press, 1993.
- [11] E. Hunt, M. P. Atkinson and R. W. Irving, "Database indexing for large DNA and protein sequence collections," *The VLDB Journal*, Vol. 11, No. 3, pp. 256-271, 2002.
- [12] T. Kahveci and A. K. Singh, "An Efficient Index Structure for String Databases," In Proceedings of the 27th VLDB Conference, pp. 351-360, 2001.
- [13] K. Kelly and P. Labute, "The A* Search and Applications to Sequence Alignment," <http://www.chemcomp.com/article/astar.htm>, 1996.
- [14] S. Kurtz, Reducing the Space Requirement of Suffix Trees. *Softw. Pract. Exp.*, Vol 29, pp. 1149-1171, 1999.
- [15] S. Kurtz, J. Choudhuri, E. Ohlebusch, C. Schleiermacher, J. Stoye, and R. Giegerich, "REPuter: the manifold applications of repeat analysis on a genome scale," *Nucleic Acids Research*, Vol. 29, No. 22, pp. 4633-4642, 2001.
- [16] U. Manber and G. Myers, "Suffix Arrays: A New Method for On-Line String Searches," *SIAM J. Comput.*, Vol. 22, No. 5, pp. 935-948, 1993.
- [17] C. Meek, J. M. Patel, and S. Kasetty, "OASIS: An Online and Accurate Technique for Local-Alignment Searches on Biological sequences," In Proceedings of the 29th VLDB Conference, pp. 920-921, 2003.
- [18] G. Navarro and R. Baeza-Yates, "A Hybrid Indexing Method for Approximate String Matching," *J. of Discrete Algorithms*, Vol. 1, No. 1, pp.205-239, 2000.
- [19] <http://www.ncbi.nlm.nih.gov>
- [20] K. Sadakane and T. Shibuya, "Indexing huge genome sequences for solving various problems," In Proceedings of the 12th

Genome Informatics, pp. 175-183, 2001.

[21] H. Shang and T. H. Merrett, "Tries for approximate string matching," IEEE Trans. on Knowledge and Data Engineering, Vol. 8, No. 4, pp. 540-547, 1996.

[22] T. Smith and M. Waterman, "Identification of Common Molecular Subsequences," Journal of Molecular Biology, 147, pp. 195-197, 1981.

[23] G. A. Stephen, String Searching Algorithms, World Scientific Publishing, 1994.

[24] Z. Tan, X. Cao, B. Ooi, and A. Tung, "The ed-tree: An Index for Large DNA Sequence Databases," In Proceedings of SSDBM Conference, pp. 1-10, 2003.

[25] E. Ukkonen, "Approximate string matching over suffix trees," In Proceedings of Combinatorial Pattern Matching (CPM93), pp. 228-242, 1993.

[26] H. Wang et al., "BLAST++: A Tool for BLASTing Queries in Batches," In Proceedings First Asia-Pacific Bioinformatics Conference, pp. 71-79, 2003.

[27] H. E. Williams and J. Zobel, "Indexing and Retrieval for Genomic Databases," IEEE TKDE Vol. 14, No. 1. pp. 63-78, 2002.



2005. 2 :
()
2005. 3~ :
:

, XML,
e-mail: kyoons@hallym.ac.kr



1992. 2 :
()
1997. 8 :
()
2003. 2 :

()
2000. 3~2004. 2 :
2004. 3~ :
:
, XML , ,

e-mail: jiwon@cs.yonsei.ac.kr



1982. 2 :
()
1985. 3 :
()
1988 3 :
()

1998. 3~1999. 2 : UCLA

1988. 4~ :

: ,
, XML, /GIS
e-mail: jhyoon@hallym.ac.kr



1989. 2 :
()
1991. 2 :
()
2001. 2 : UCLA

()
1991. 3~1996. 8 :
2001. 2~2002. 6 : IBM T. J. Watson Research
Center Post-Doctoral Fellow
2002. 8~2003. 8 :

2003. 9~ :

: , ,
, XML
e-mail: sanghyun@cs.yonsei.ac.kr



1989. 2 :
()
1991. 2 :
()
1994. 2 :

()
1991. 7~8 : Stanford University, Computer

Science Department

1994. 2~1995. 2 : KAIST

1999. 8~2000. 8 : IBM T.J. Watson
Research Center Post-Doc.

1995. 3~2000. 8 :

2003. 3~ :

: , ,
, , , /GIS,
,
e-mail: wook@hanyang.ac.kr

