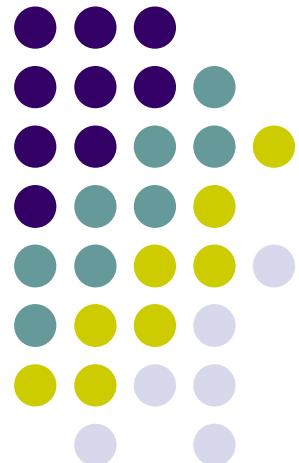


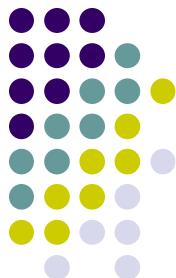
Directed Acyclic Word Graphs and Their Efficient Implementations

Shunsuke Inenaga

*Department of Computer Science
University of Helsinki*



Short CV



Shunsuke Inenaga

March 2003

***PhD Kyushu University, Japan
“String Processing Algorithms”***

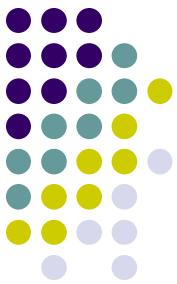
April – Sept. 2003

Posdoc of Japan Science Technology Agency (JST)

Sept. 2003 -

***Posdoc of Dep. of Computer Science,
University of Helsinki***

Pattern Matching Problem



Input: text string T and pattern string P .

Output: all occurrences of P in T .

When T is fixed

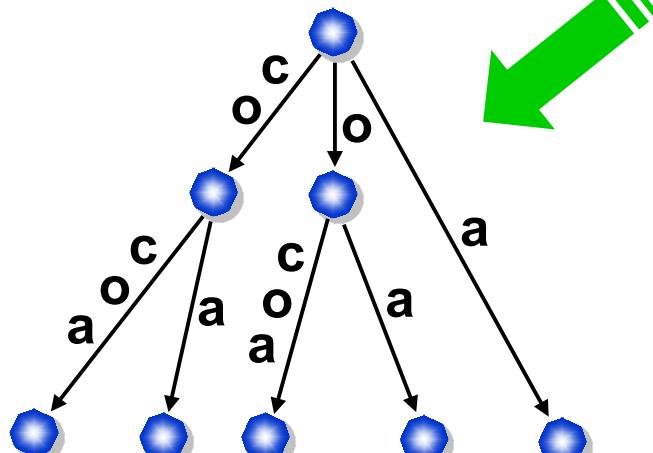
Text Indexing Structures

Text Indexing Structures

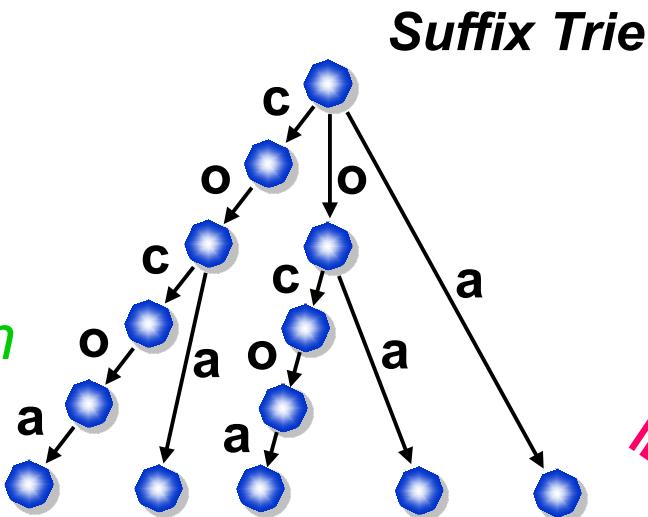


Text: **cocoa**

Suffix Tree

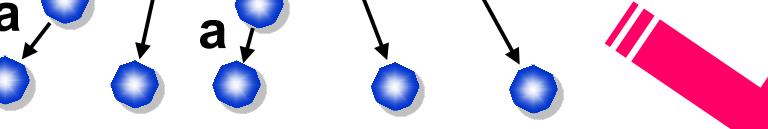


Compaction

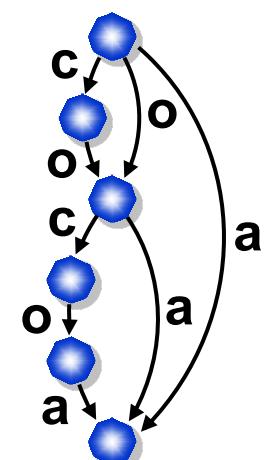


Suffix Trie

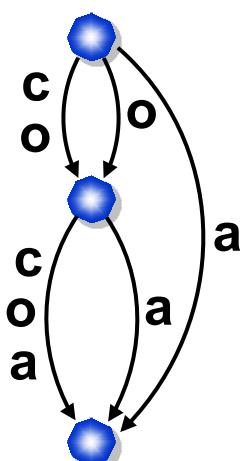
Directed Acyclic Word Graph



Minimization



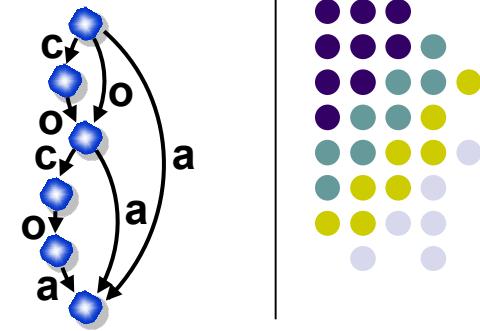
Minimization



Compaction

Compact Directed Acyclic Word Graph

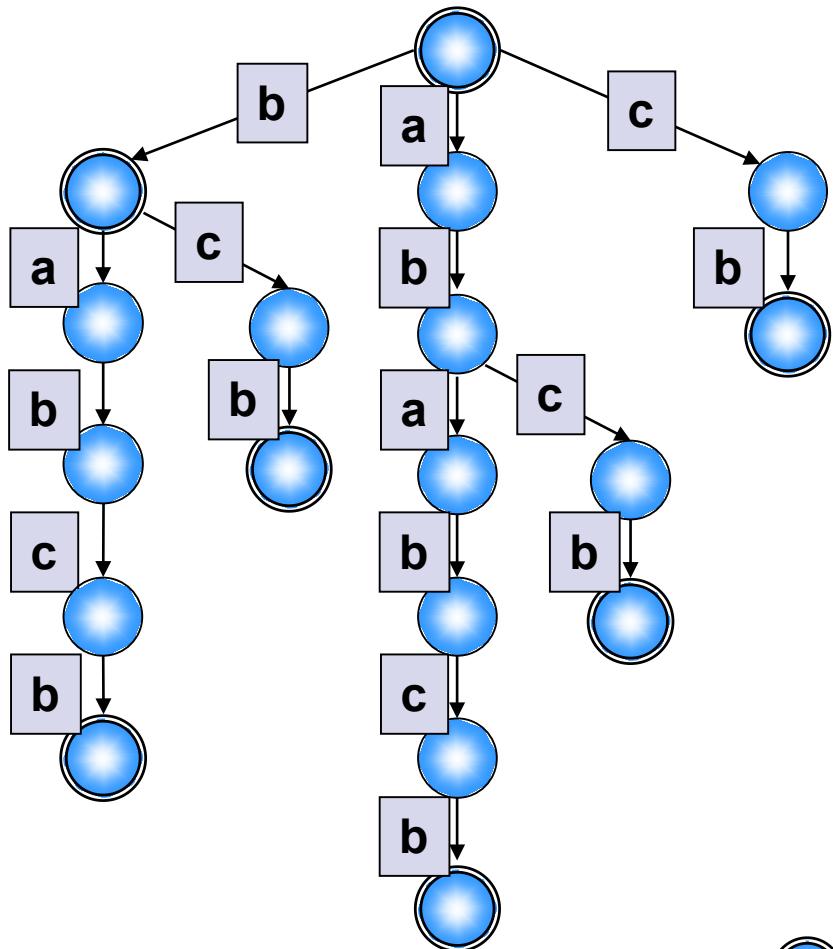
Directed Acyclic Word Graphs



The Directed Acyclic Word Graph (DAWG) of string w is the smallest DFA that accepts all suffixes of w .

DAWG = Suffix Automaton

Conversion of Suffix Trie to DAWG



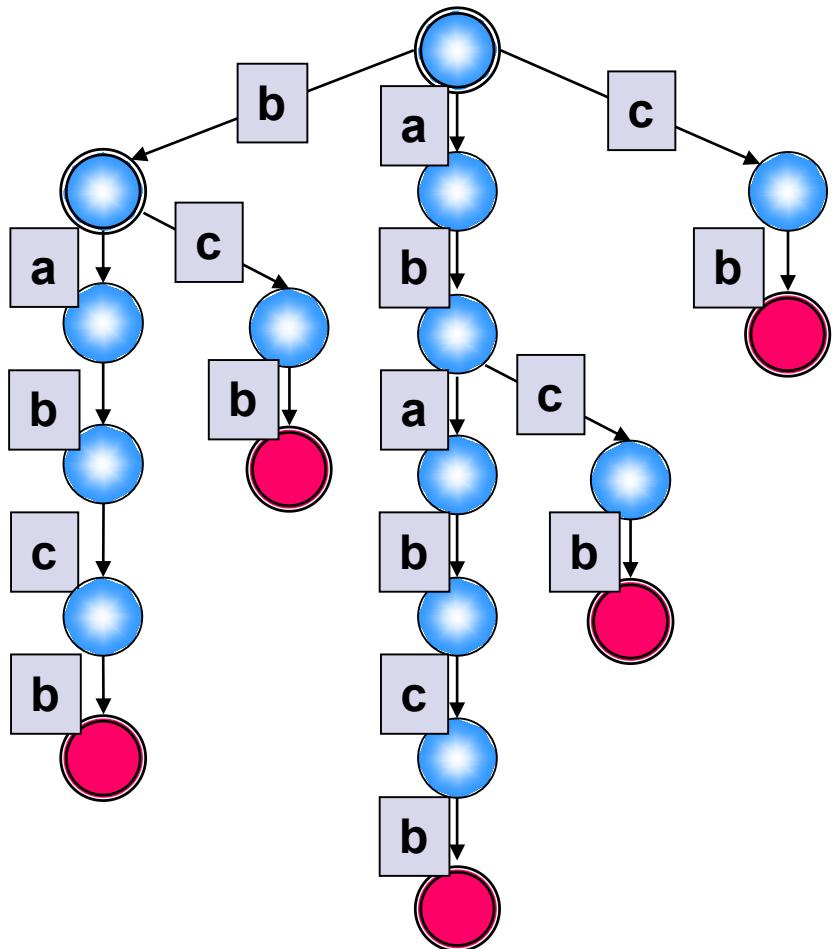
ababcb

*Merge isomorphic
subtrees*



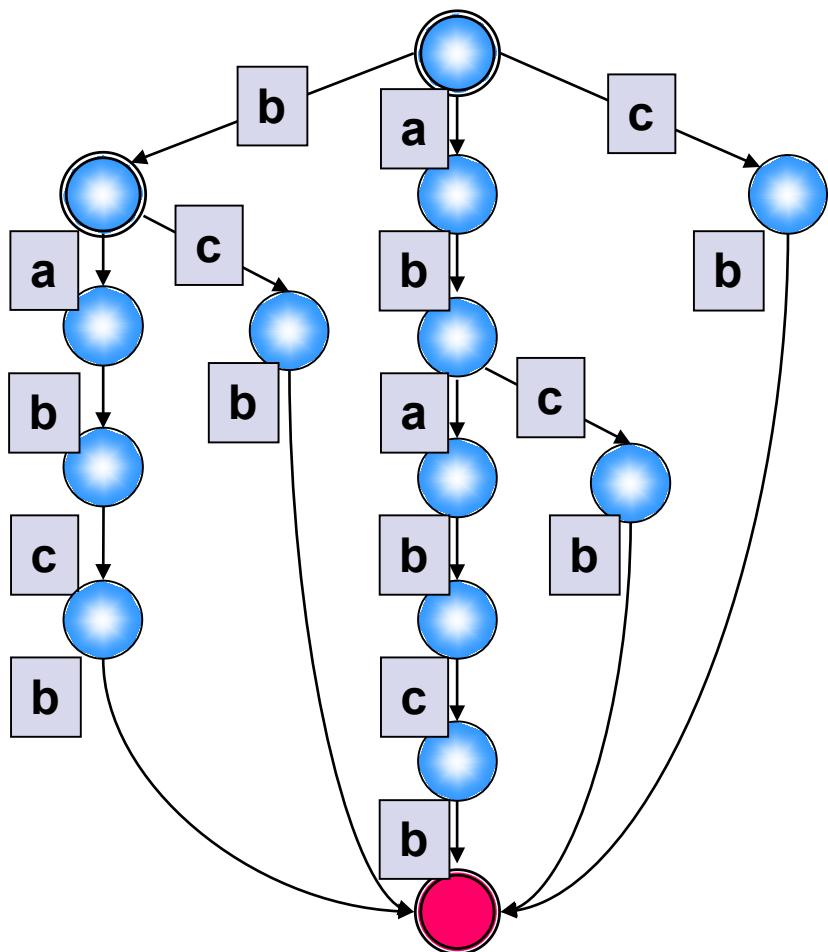
accepting node for suffixes

Conversion to DAWG



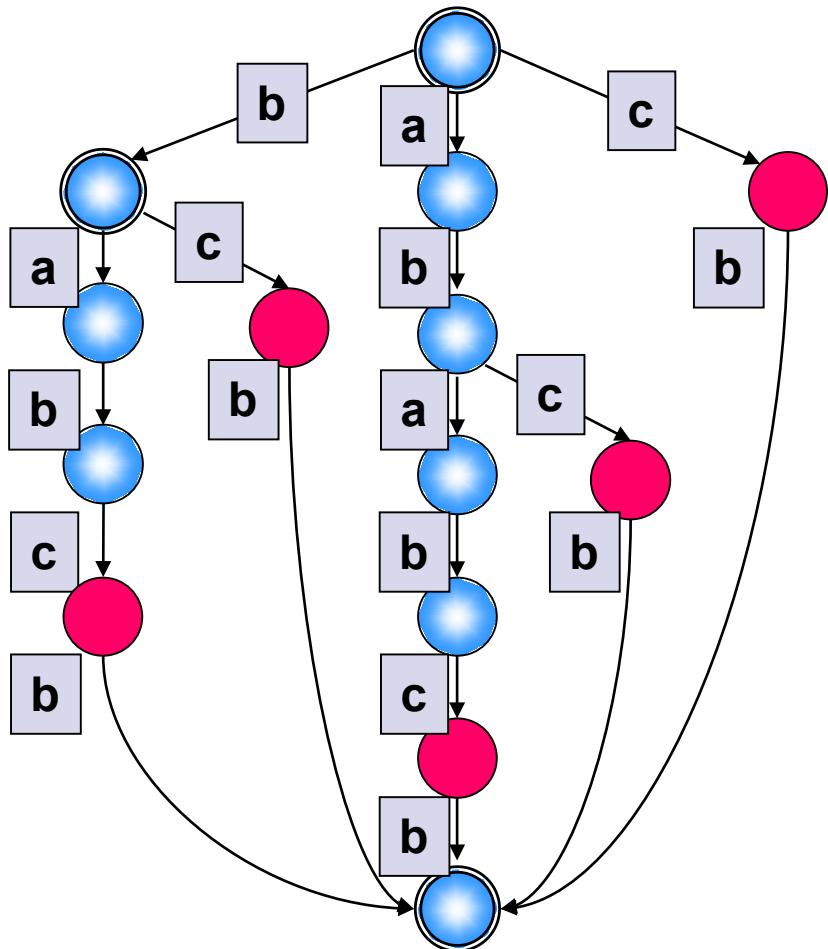
ababcb

Conversion to DAWG [Cont.]



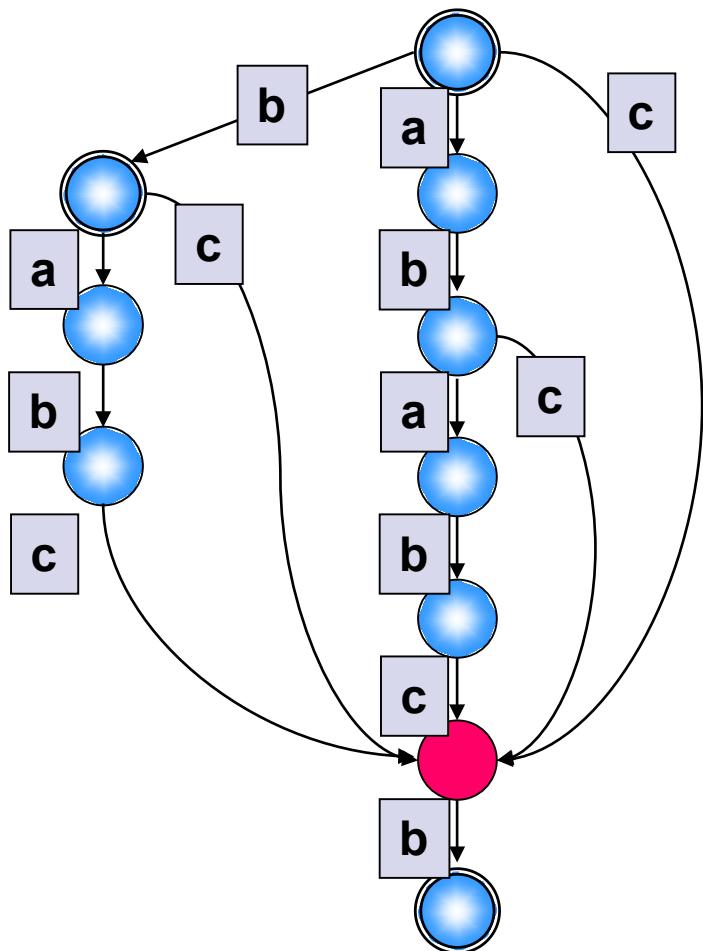
ababcb

Conversion to DAWG [Cont.]



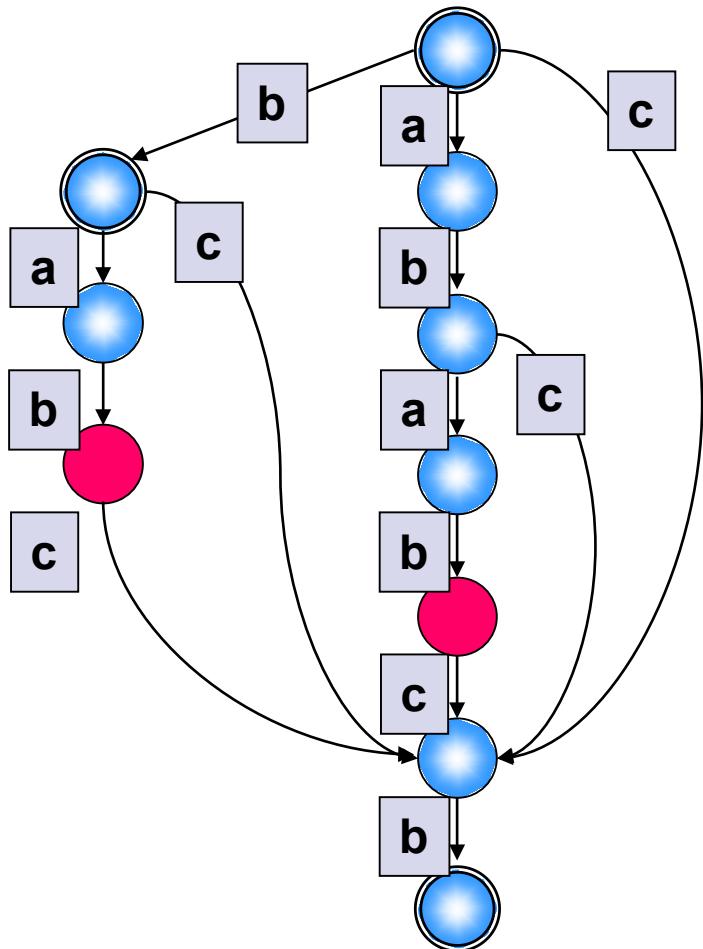
ababcb

Conversion to DAWG [Cont.]



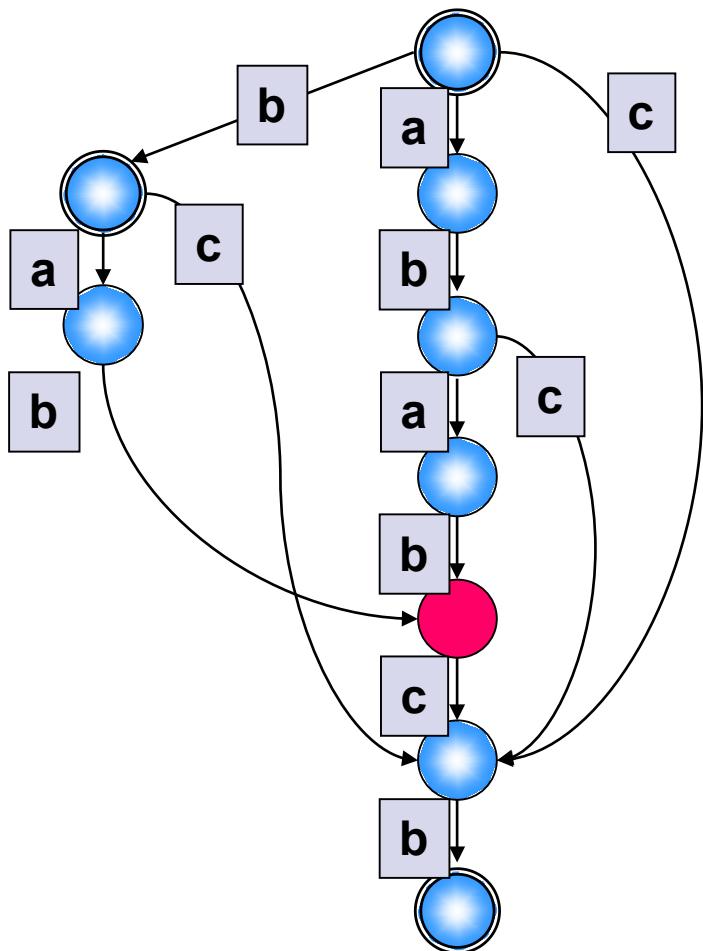
ababcb

Conversion to DAWG [Cont.]



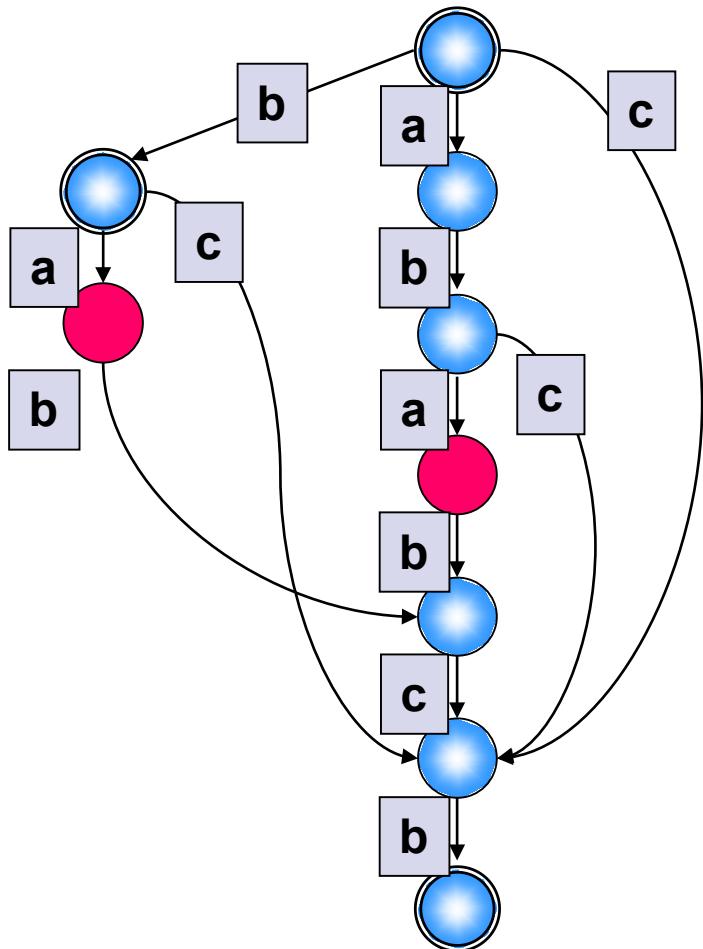
ababcb

Conversion to DAWG [Cont.]



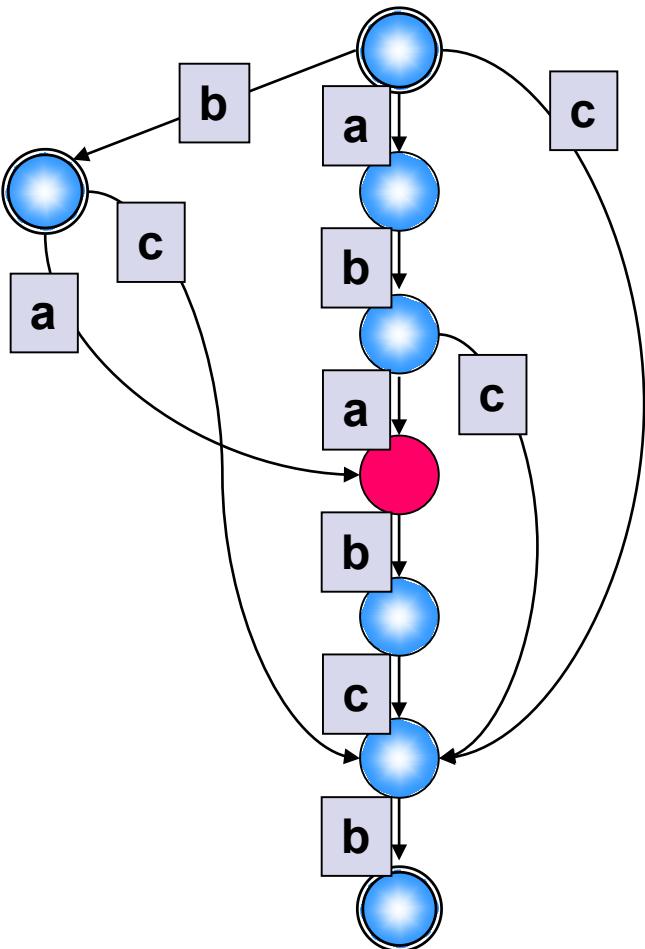
ababcb

Conversion to DAWG [Cont.]



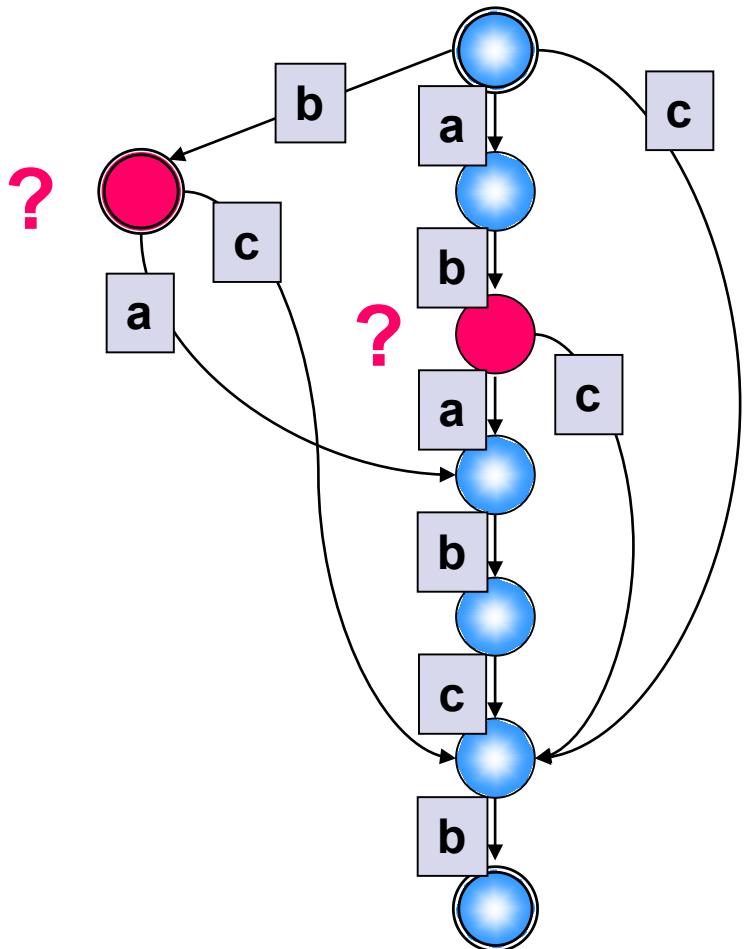
ababcb

Conversion to DAWG [Cont.]



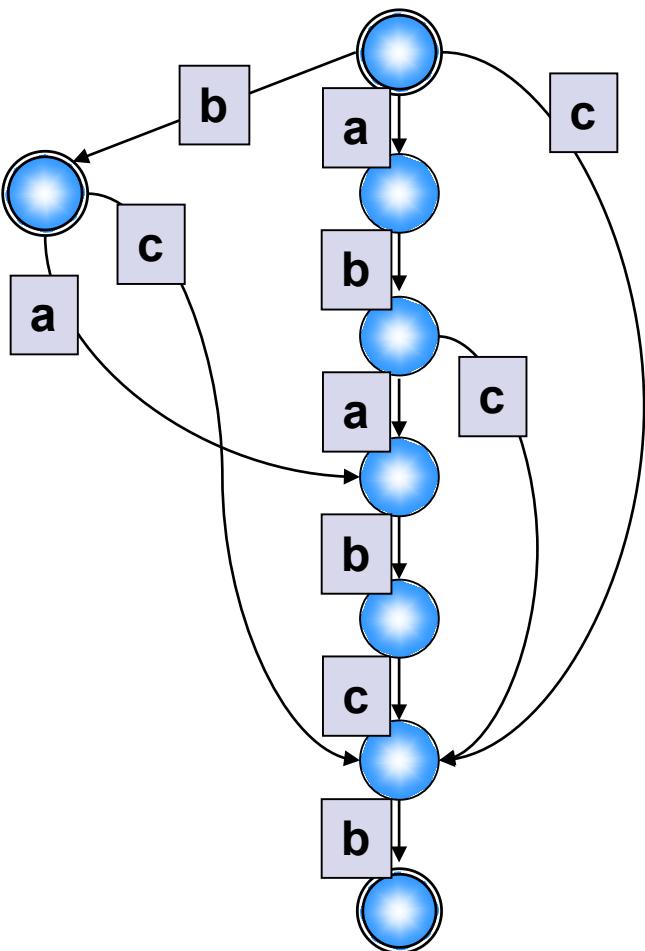
ababcb

Conversion to DAWG [Cont.]



ababcb

DAWG Completed

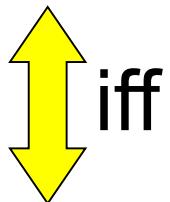


ababcb

Equivalence Class on DAWG



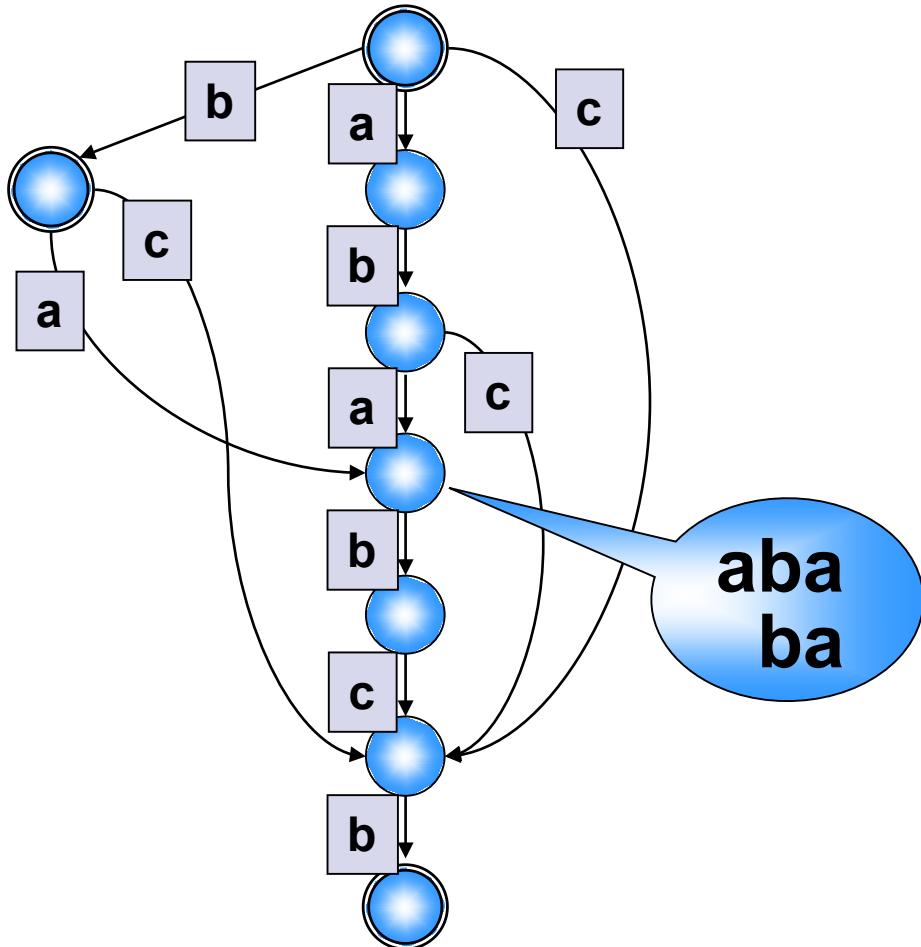
*substrings x, y of w are recognized
by the same node of $DAWG(w)$*



$$EndPos_w(x) = EndPos_w(y)$$

$EndPos_w(x)$: the set of all positions of w where x ends

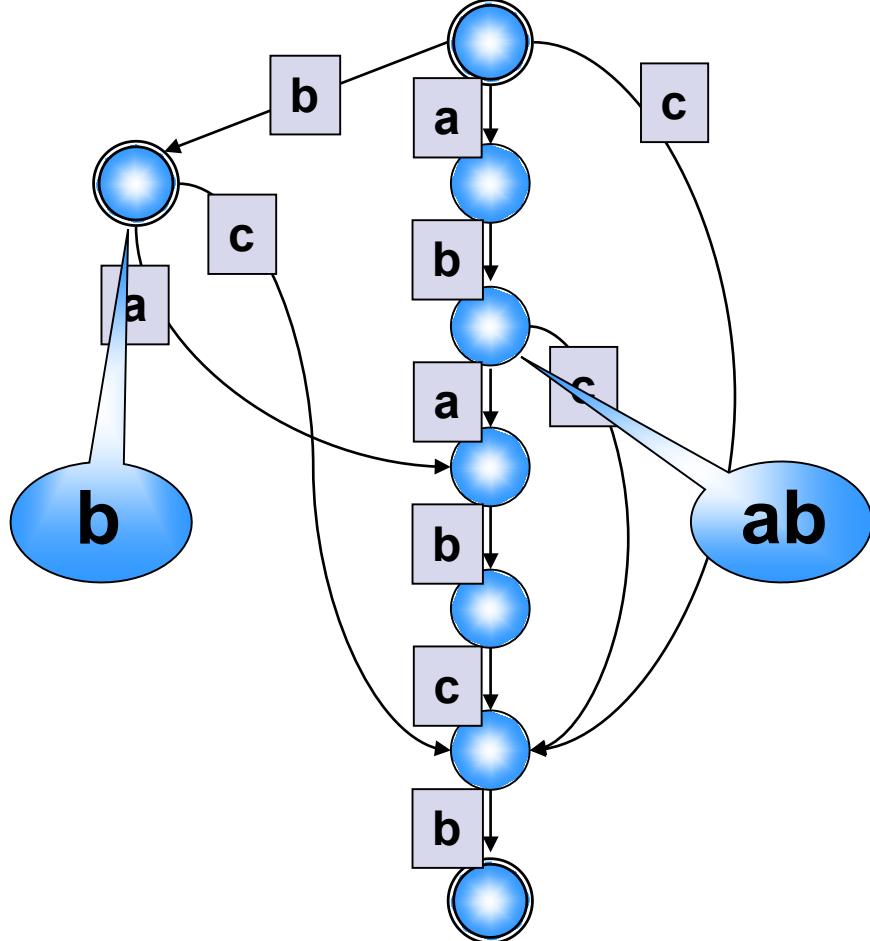
Equivalence Class on DAWG [Cont.]



ababcb
1 2 3 4 5 6

$\text{EndPos}_w(\text{aba}) = \{3\}$
 $\text{EndPos}_w(\text{ba}) = \{3\}$

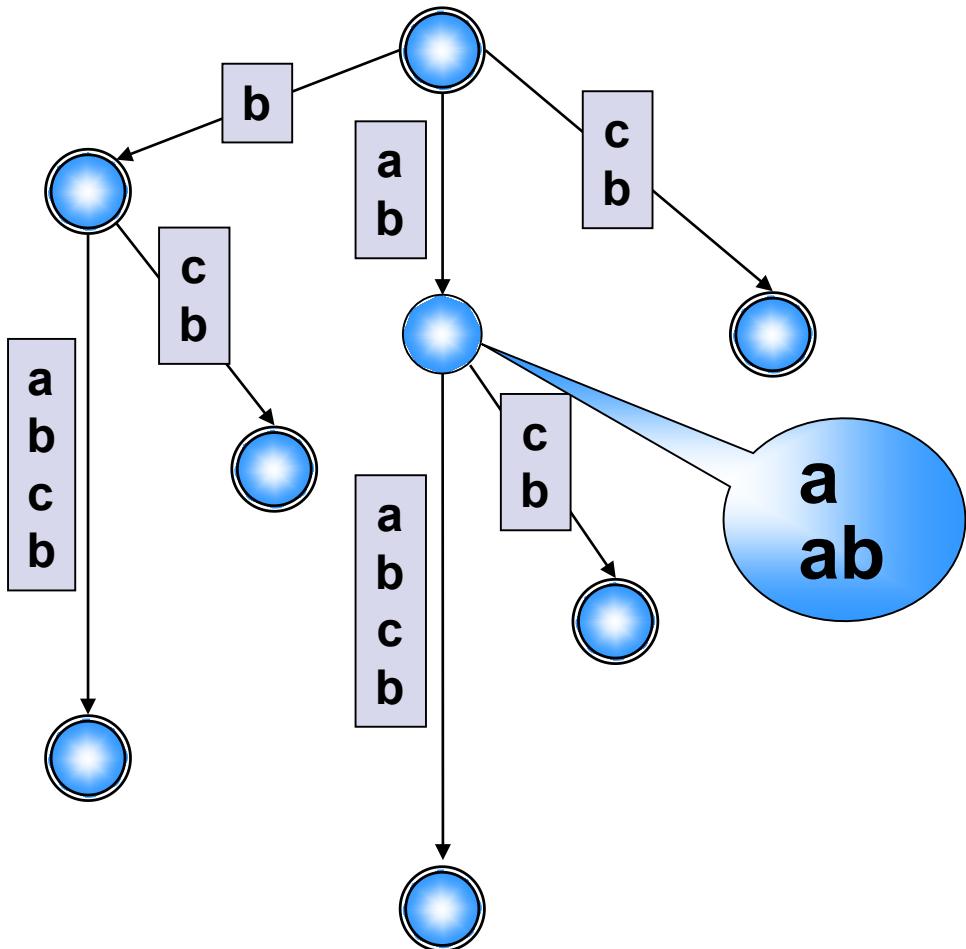
Equivalence Class on DAWG [Cont.]



ababcb
1 2 3 4 5 6

$EndPos_w(ab) = \{2,4\}$
 $EndPos_w(b) = \{2,4,6\}$

Equivalence Class on Suffix Tree



ababcb
1 2 3 4 5 6

$BegPos_w(a) = \{1,2\}$
 $BegPos_w(ab) = \{1,2\}$

Size of DAWG



Theorem. (Blumer et al. 1985)

Let $n = |w|$.

If $n > 1$, $DAWG(w)$ has at most $2n-1$ nodes and $3n-3$ edges.

If $n = 1$, $DAWG(w)$ has exactly 2 nodes and 1 edge.

Construction of DAWG



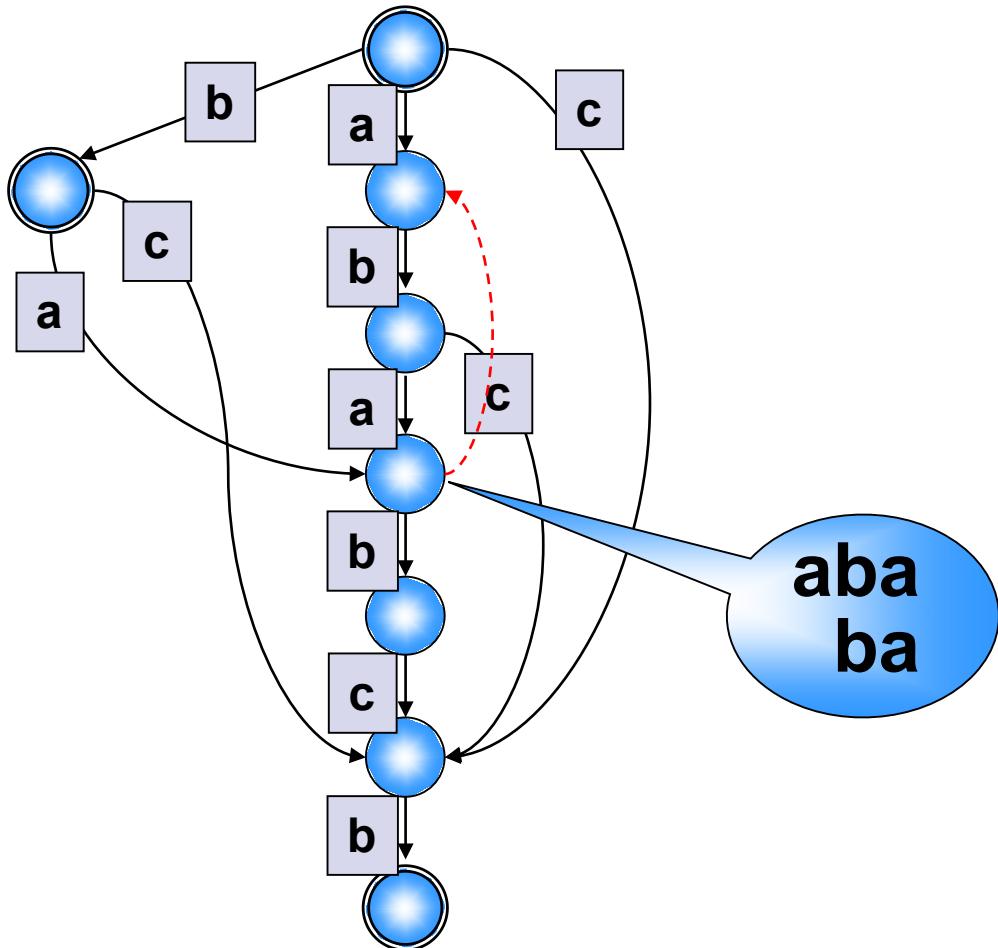
Theorem. (Blumer et al. 1985)

Assume Σ is fixed.

For any string w of length n ,

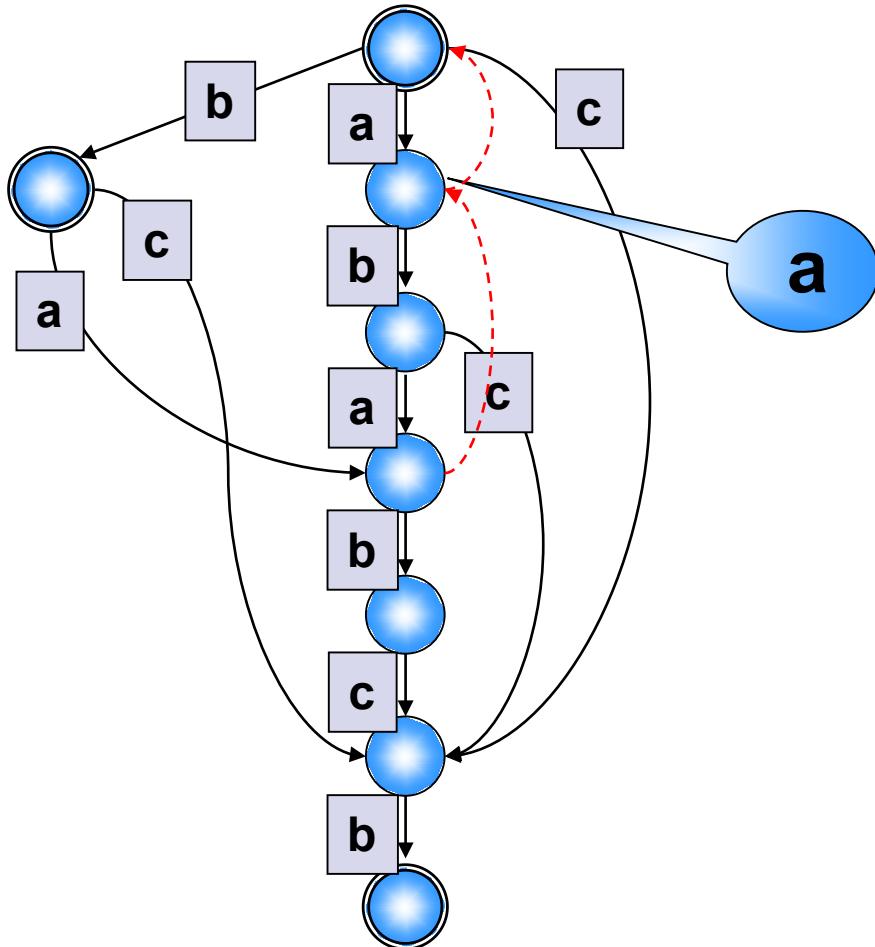
*$DAWG(w)$ can be constructed **on-line**
and in $O(n)$ time and space.*

Suffix Links



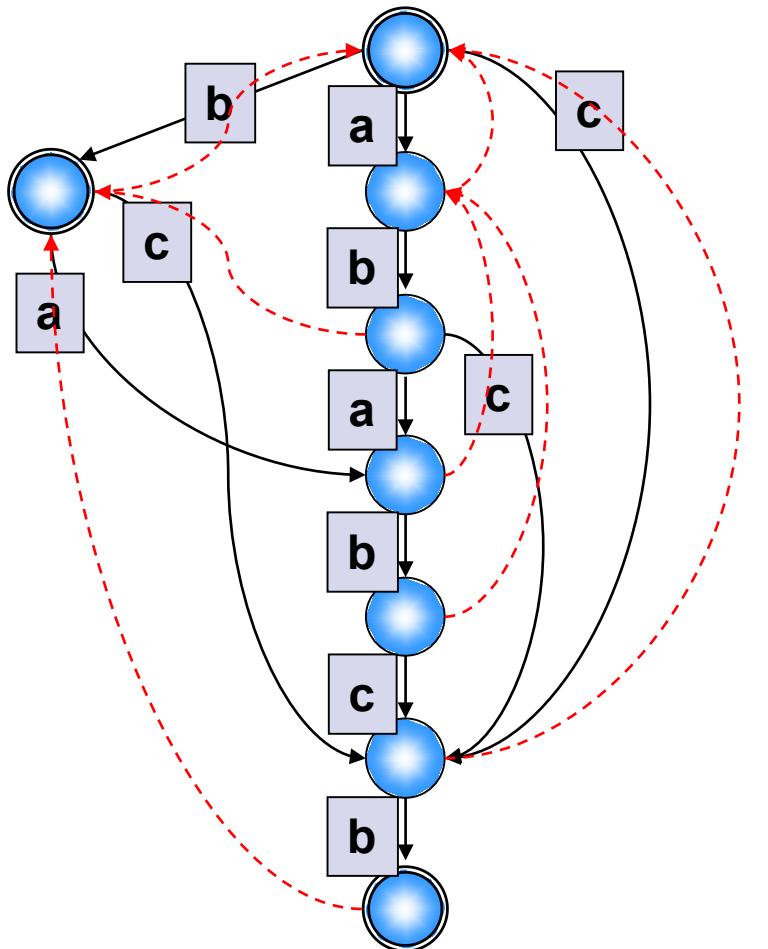
ababcb

Suffix Links



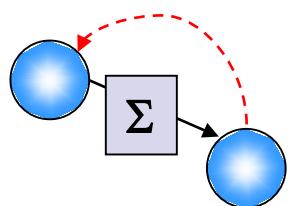
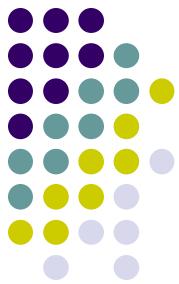
ababcb

Suffix Links



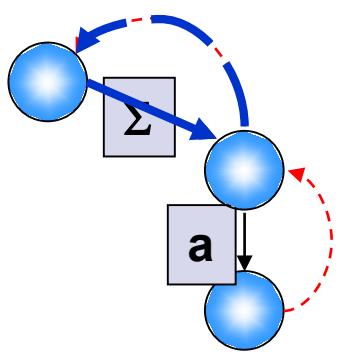
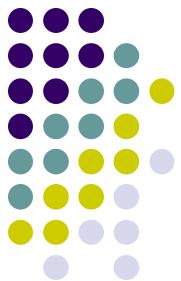
ababcb

On-line Construction of DAWG



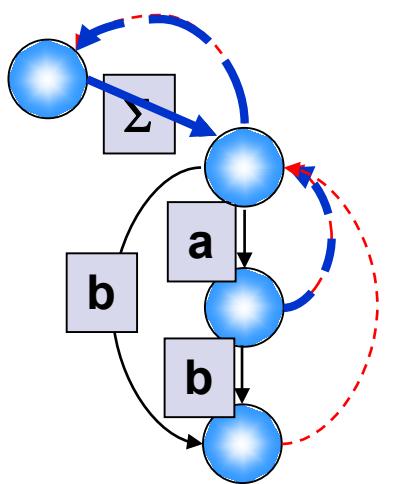
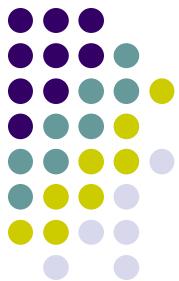
ababcb

On-line Construction of DAWG



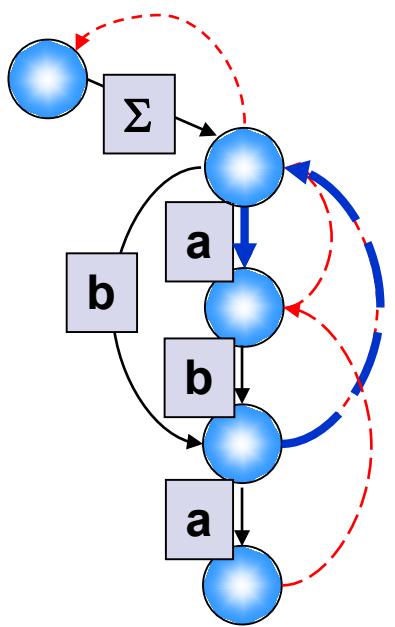
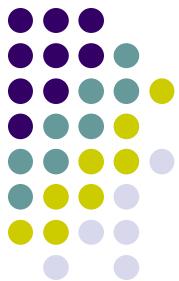
ababcb

On-line Construction of DAWG



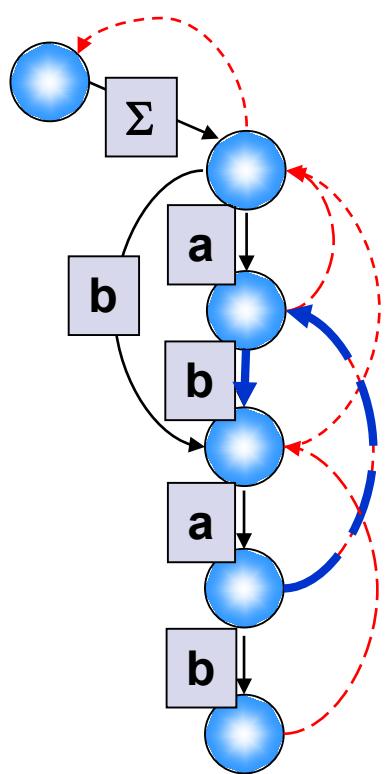
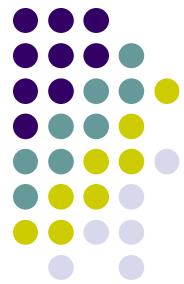
ababcb

On-line Construction of DAWG



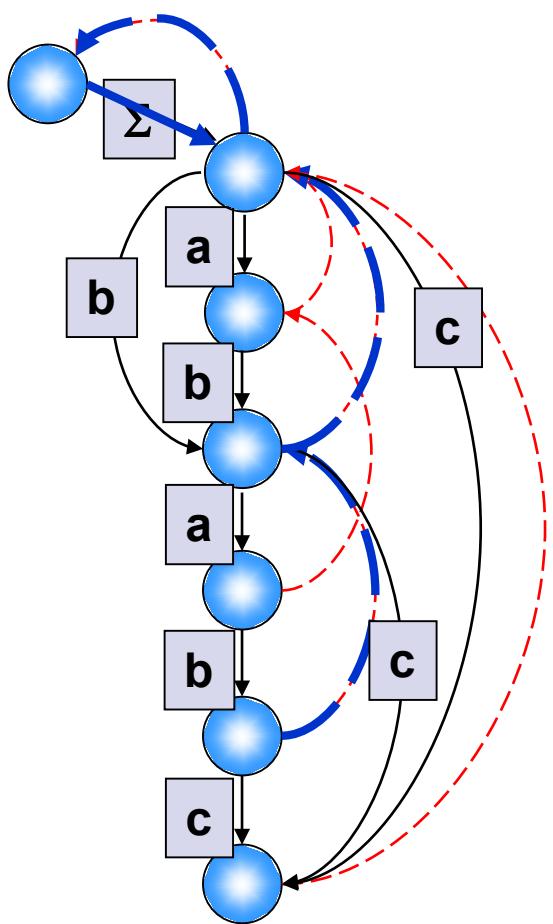
ababcb

On-line Construction of DAWG



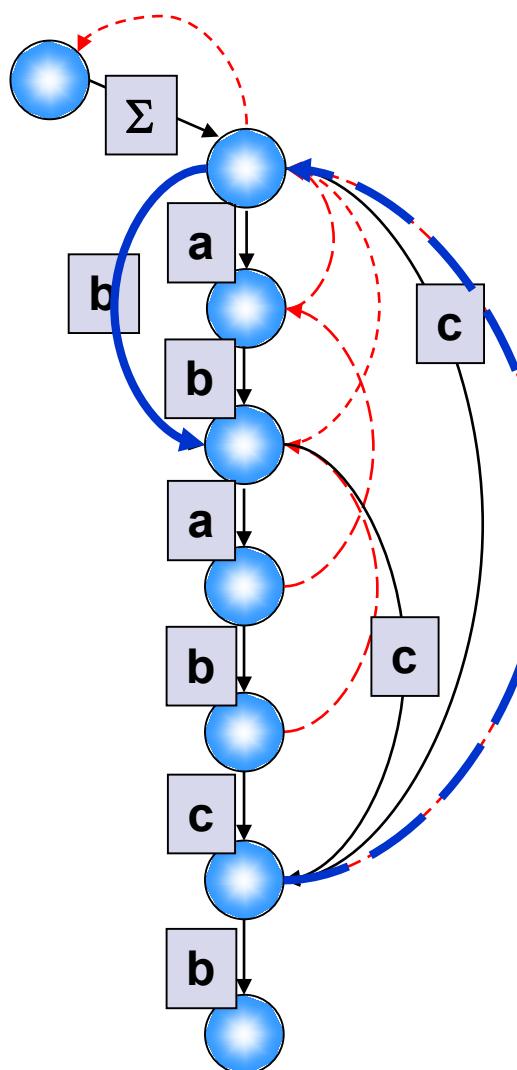
abab**c**b

On-line Construction of DAWG



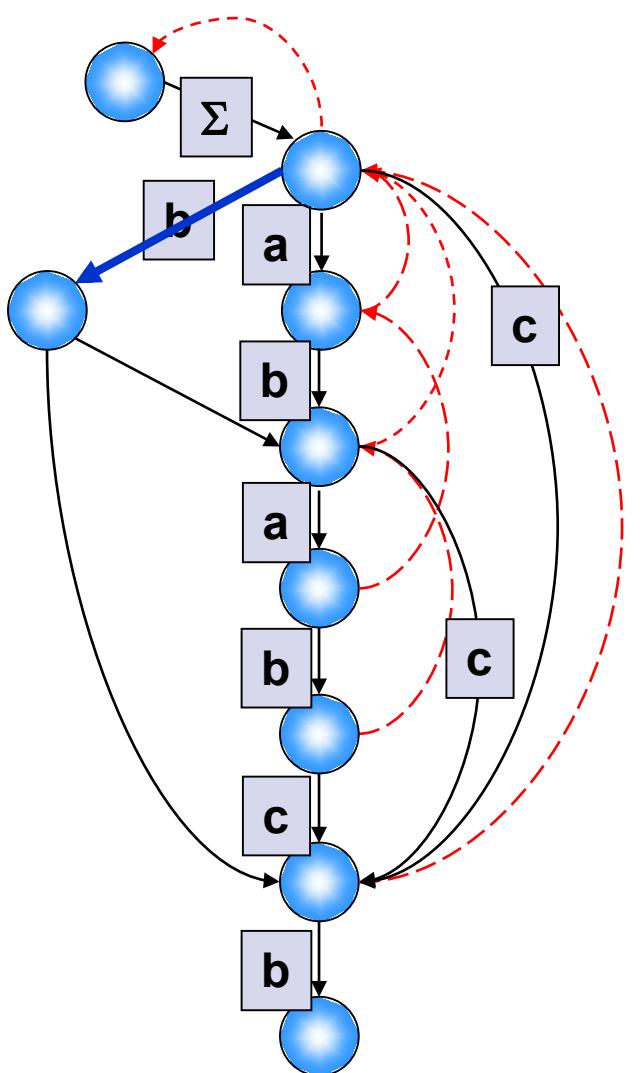
ababcb

On-line Construction of DAWG



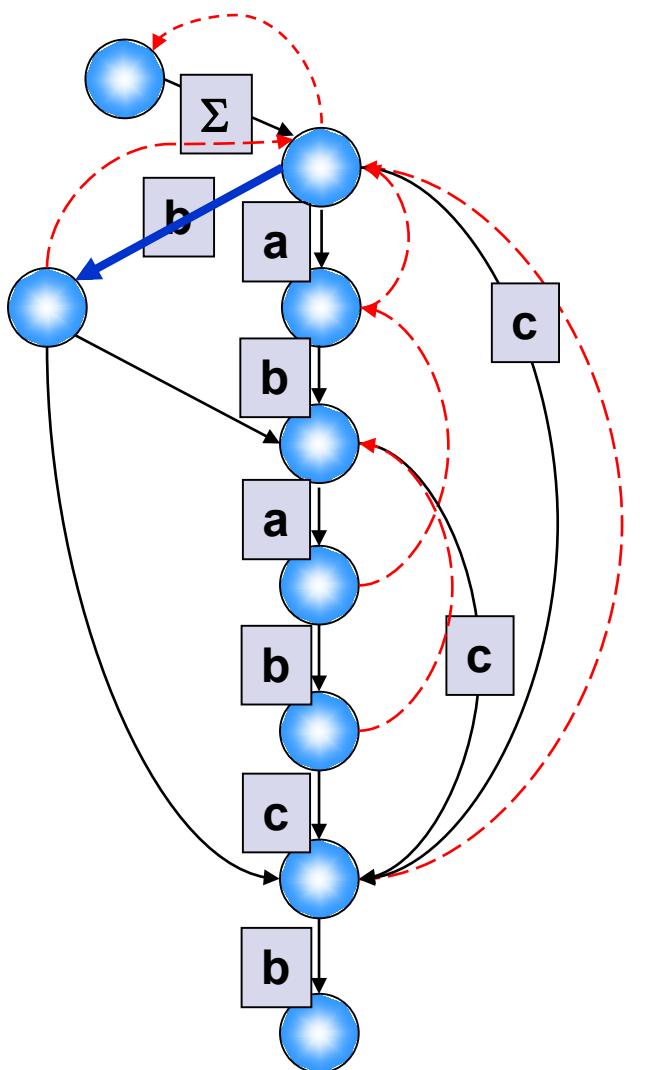
ababcb

On-line Construction of DAWG



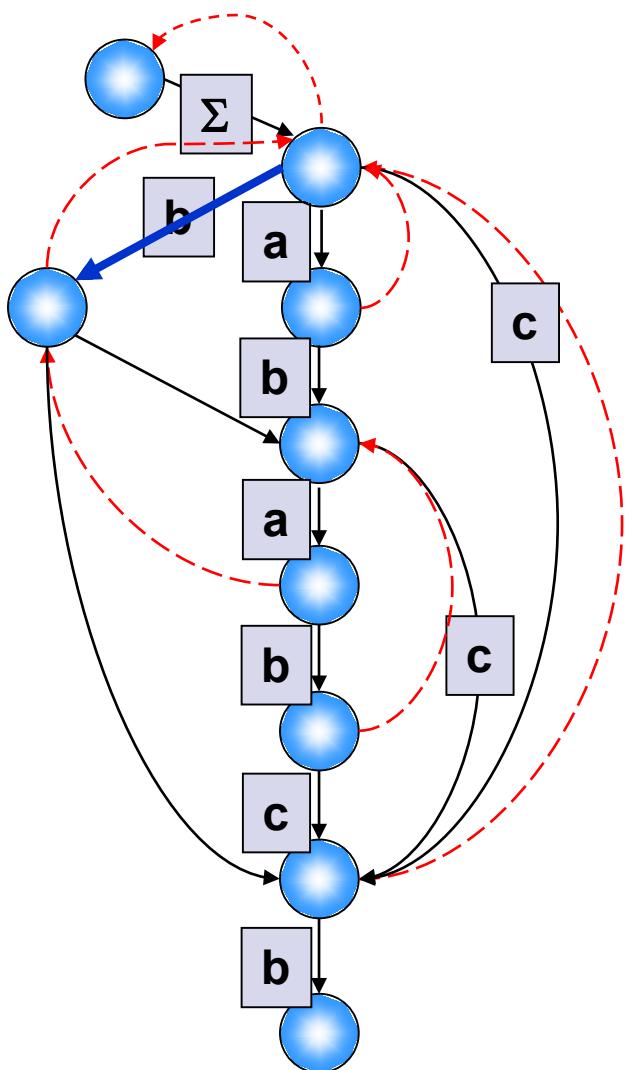
ababcb

On-line Construction of DAWG



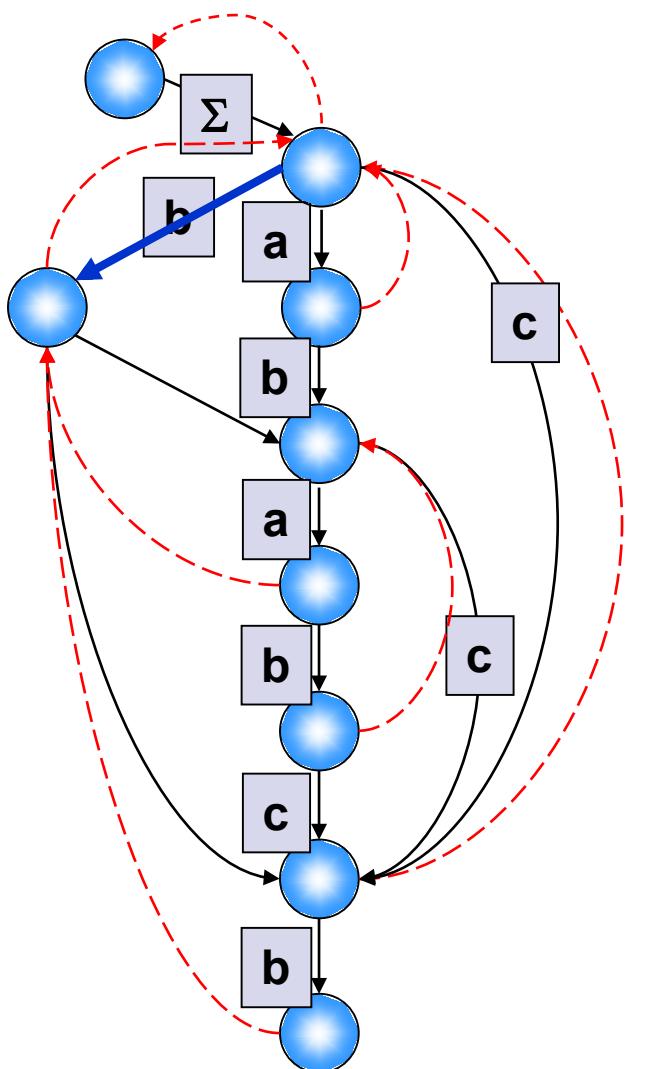
ababcb

On-line Construction of DAWG



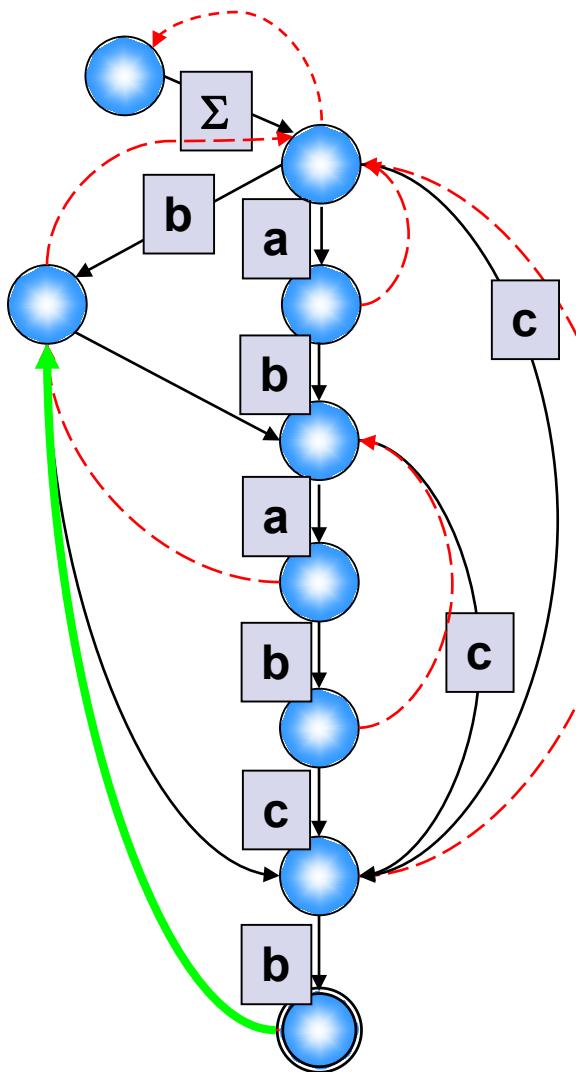
ababcb

On-line Construction of DAWG



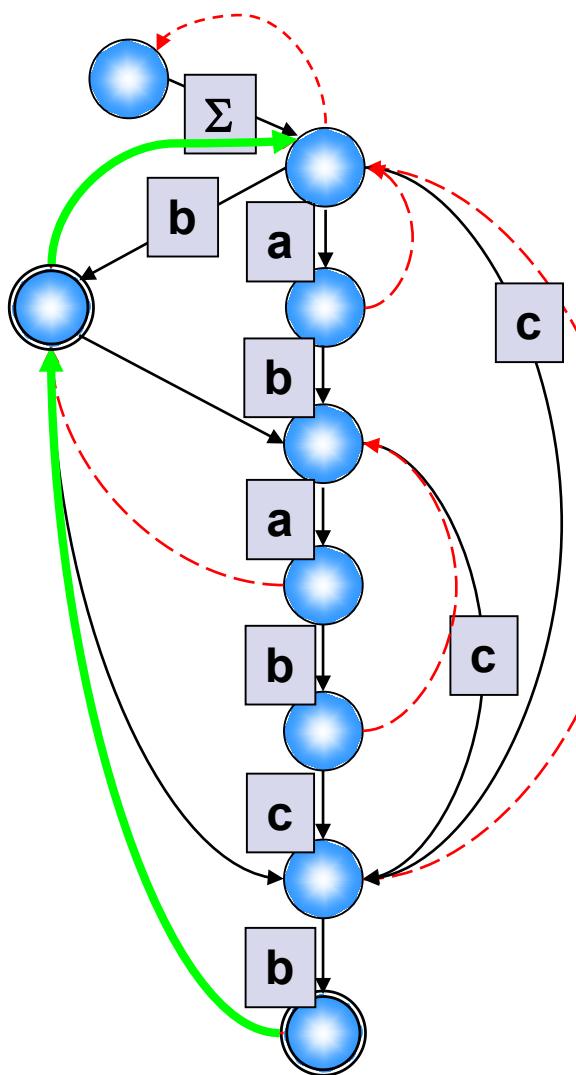
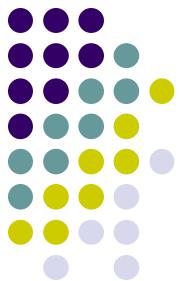
ababcb

Marking Accepting Nodes



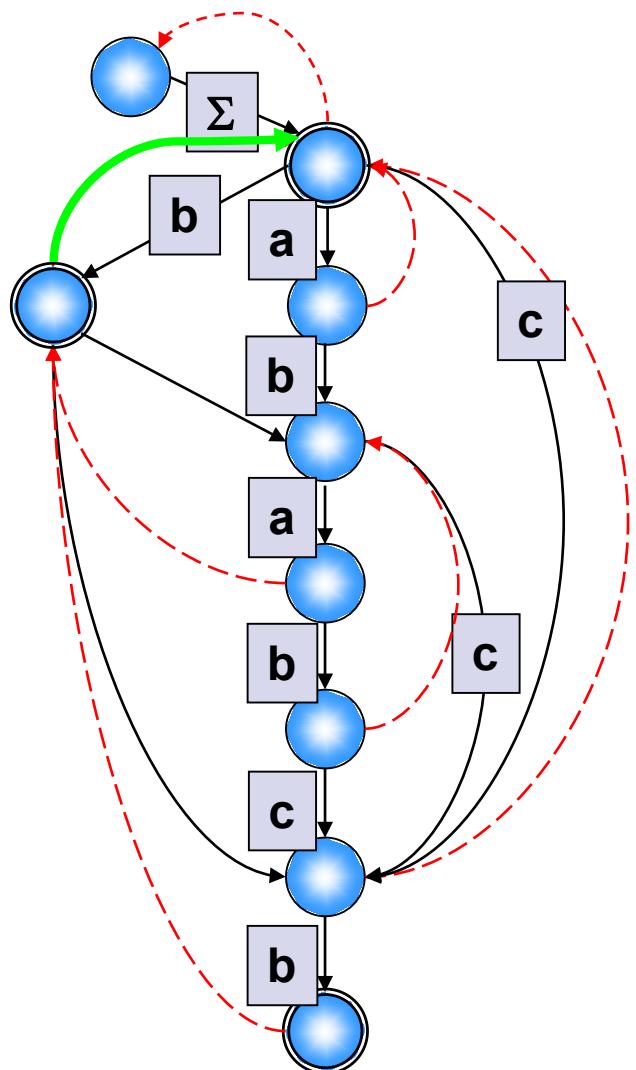
ababcb

Marking Accepting Nodes



ababcb

Marking Accepting Nodes



ababcb

Completed!!

Search Tree of a Set of Strings

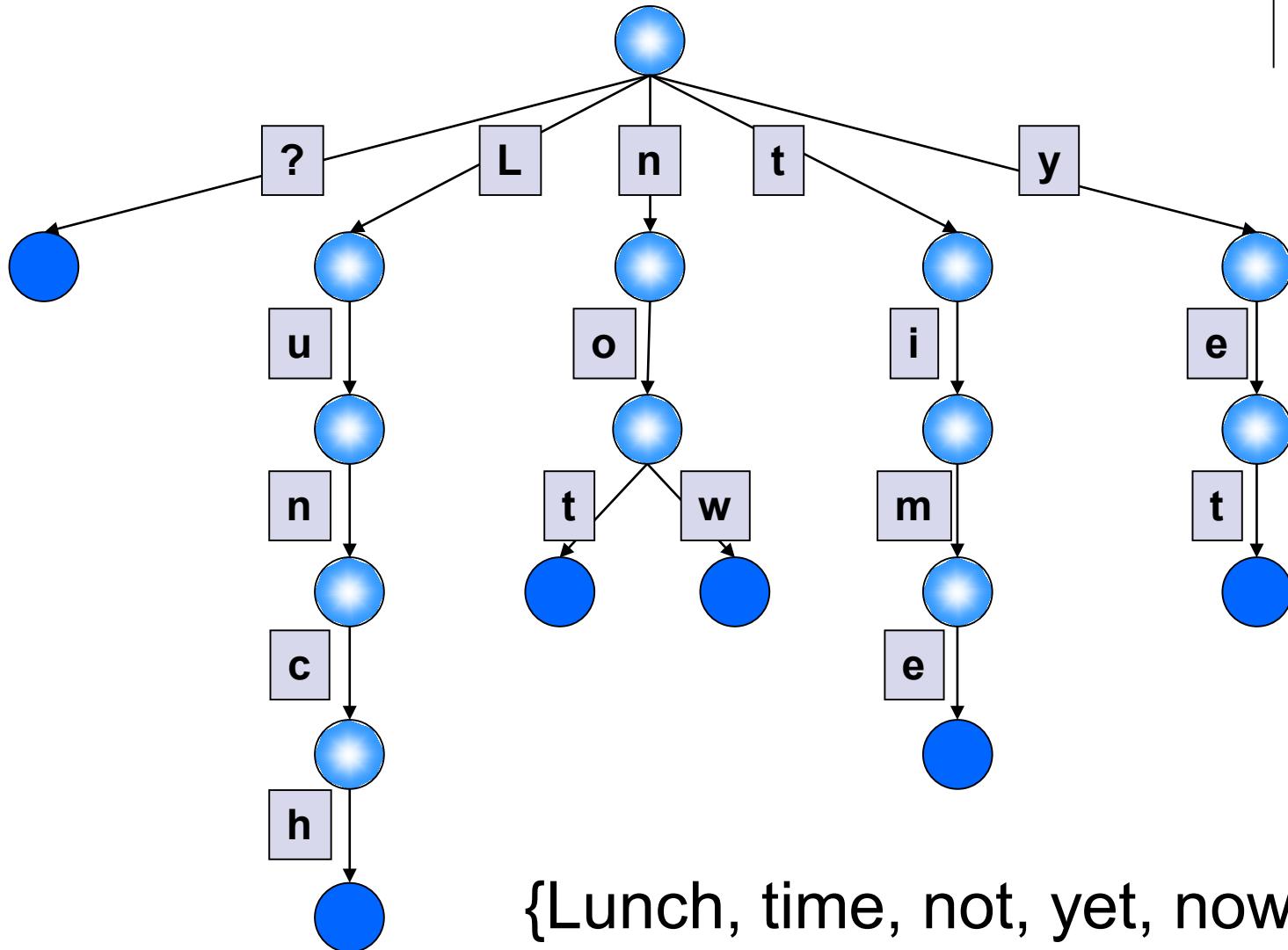


Table Implementation

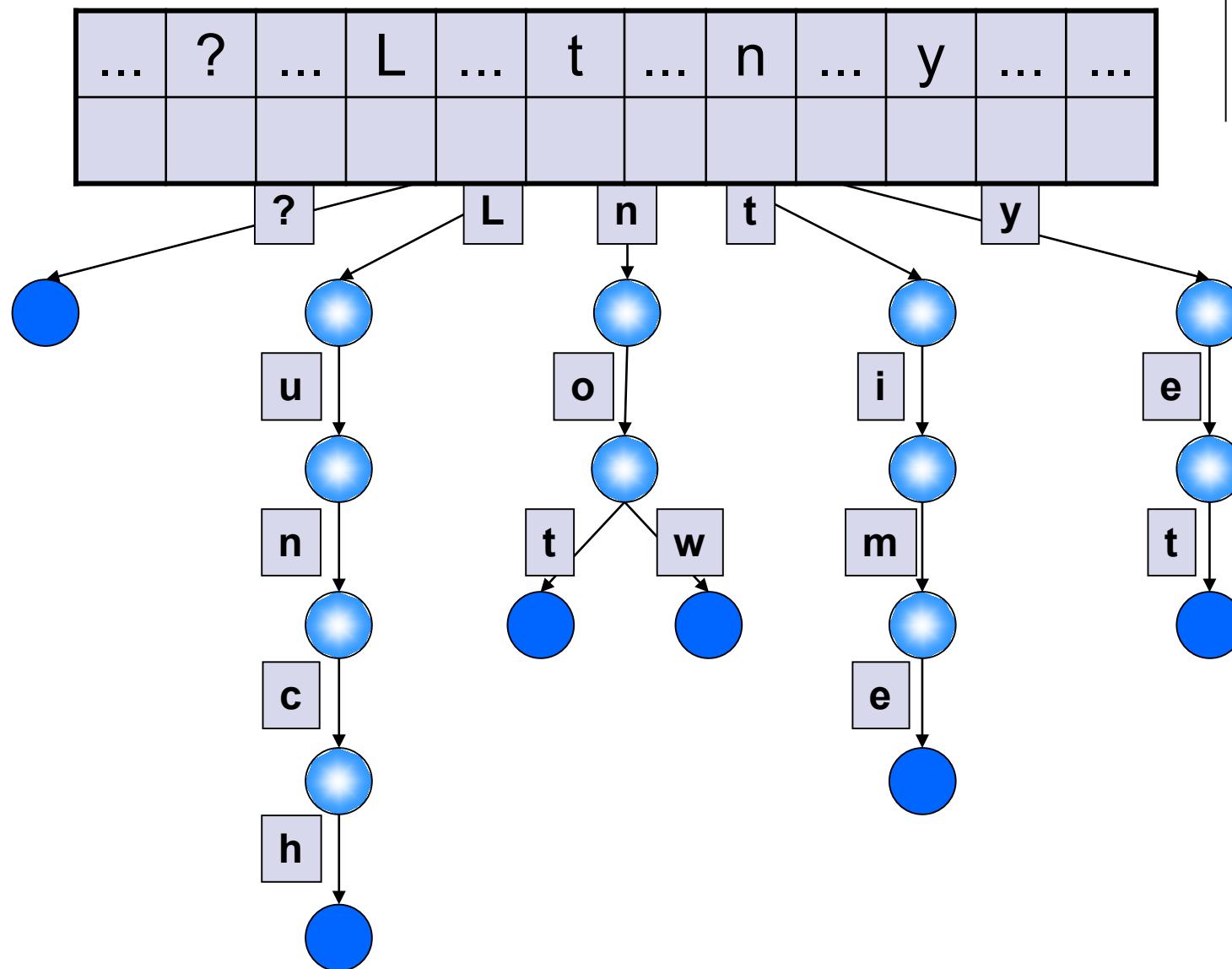
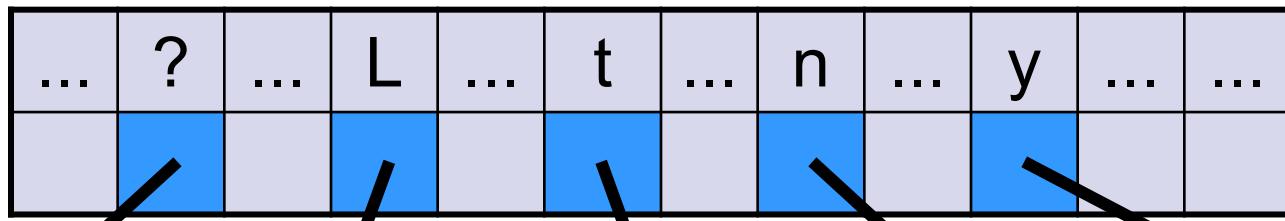
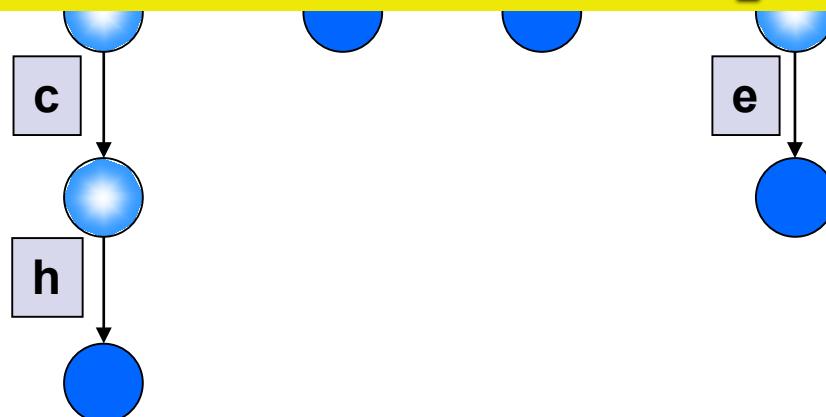


Table Implementation

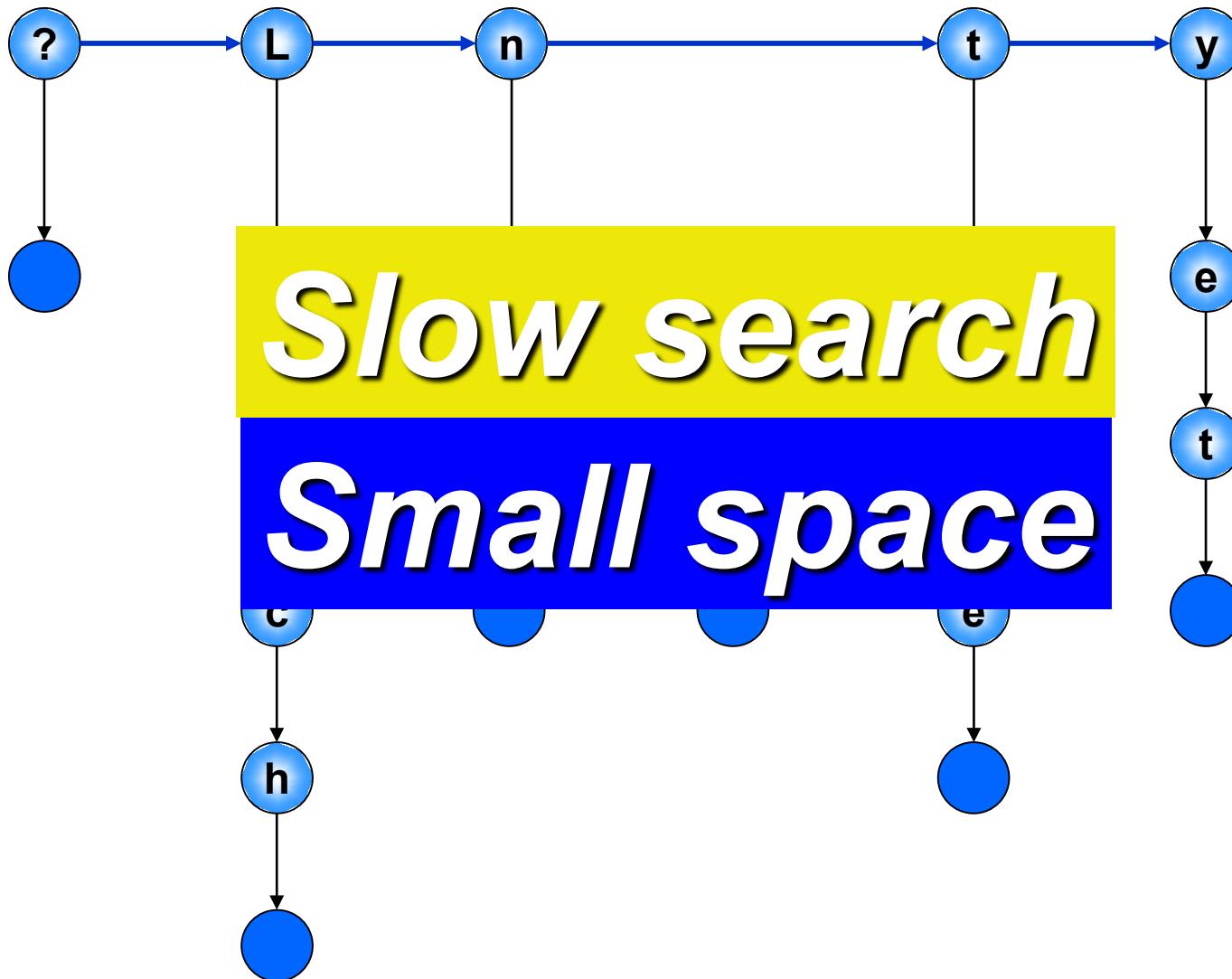


Very fast search

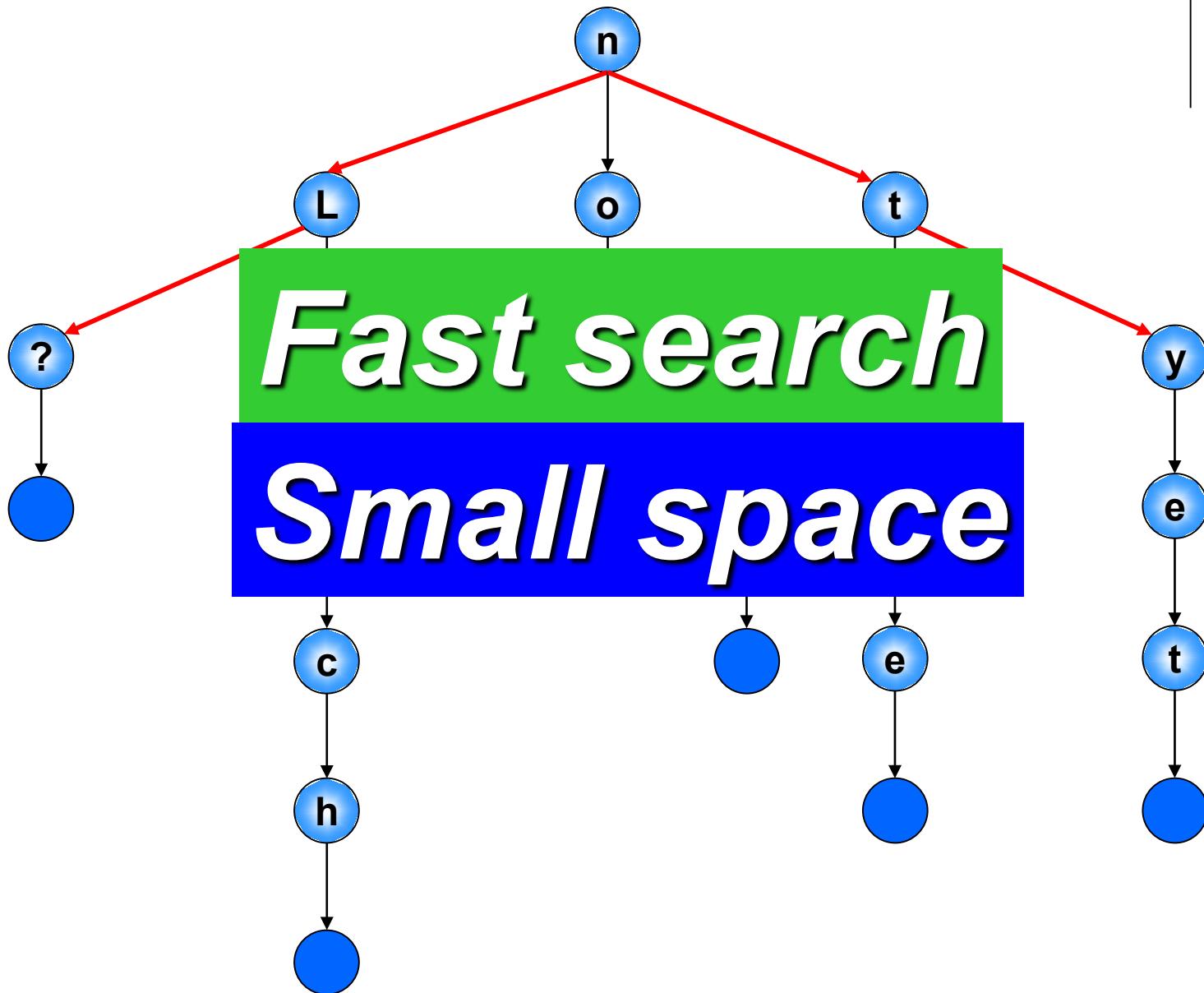
Too much space



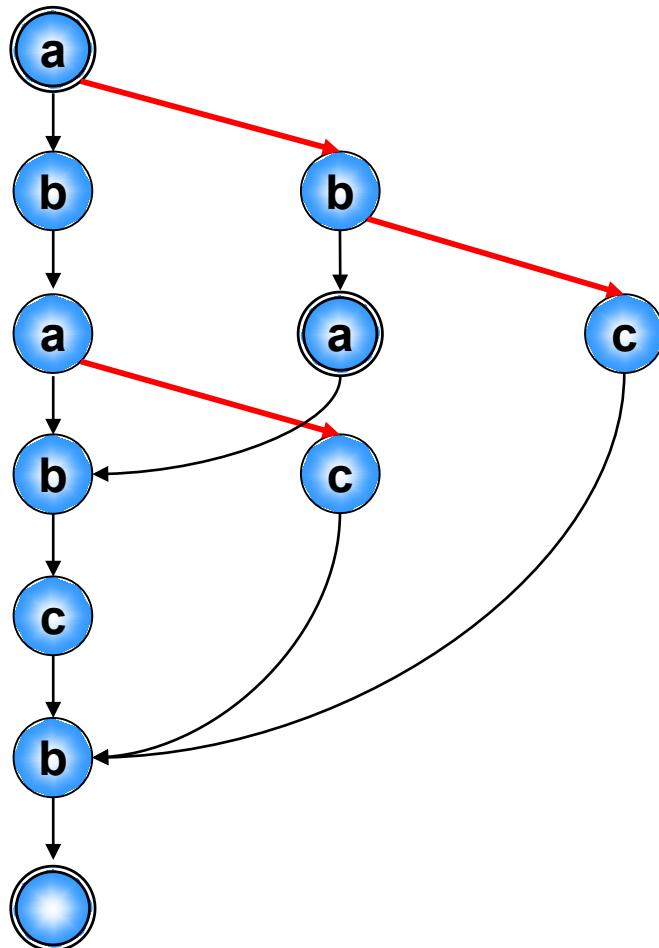
Linked List Implementation



Ternary Search Tree



Ternary Directed Acyclic Word Graph



ababcb

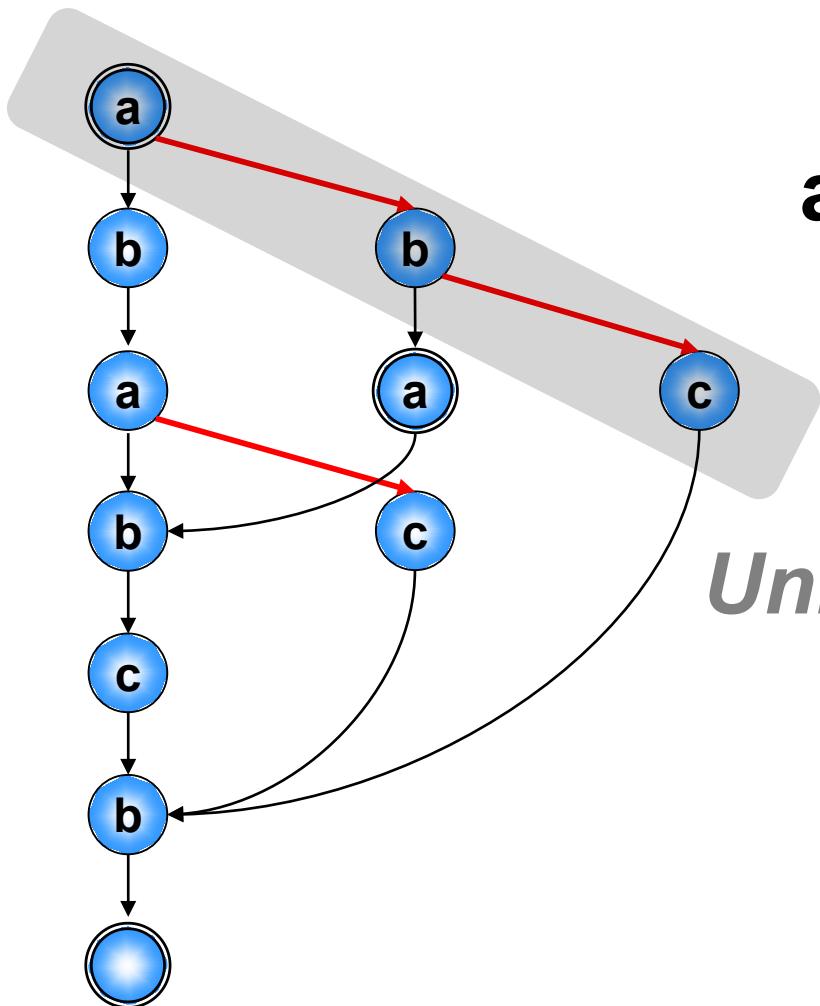
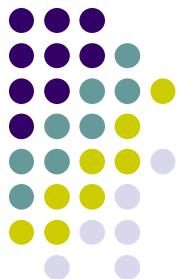
Construction of TDAWG



Theorem. (Miyamoto et al. 2003)

*For any string w of length n ,
 $TDAWG(w)$ can be constructed **on-line**
and in $O(n|\Sigma|)$ time and $O(n)$ space.*

Ternary Directed Acyclic Word Graph



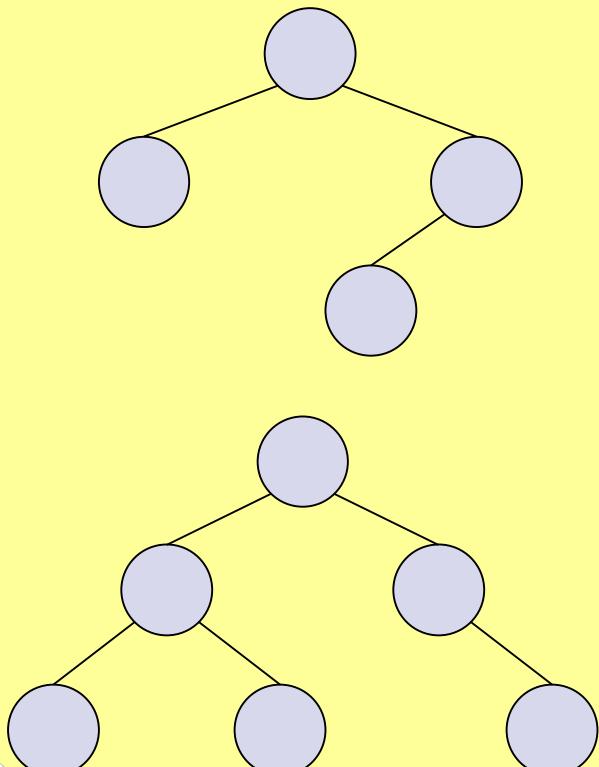
ababcb

Unbalanced!

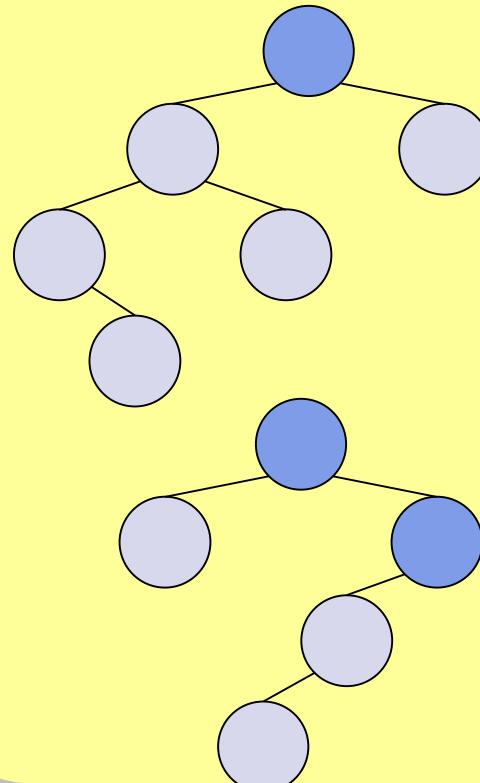
Using AVL Tree



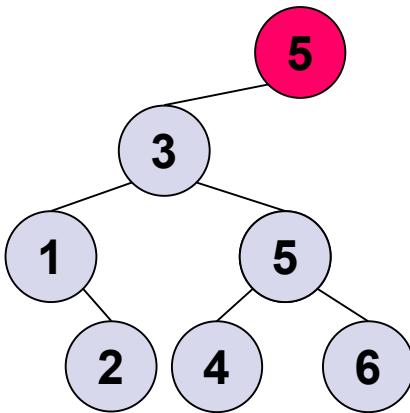
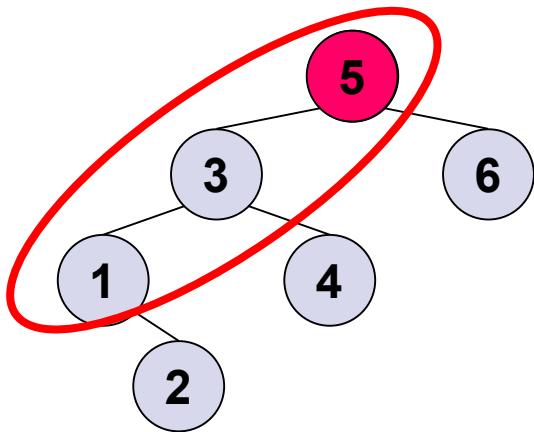
AVL Balanced



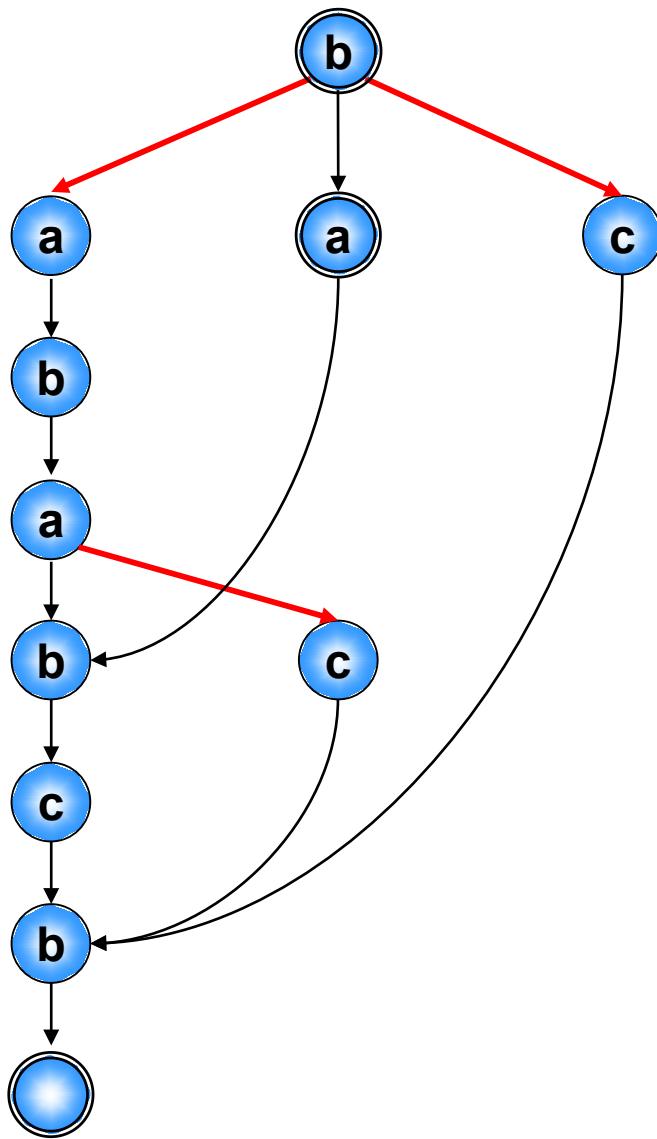
Not AVL Balanced



Node Rotation

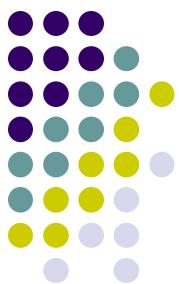


AVL DAWG



ababcb

Construction of AVL TDAWG



Theorem. (Miyamoto et al. 2003)

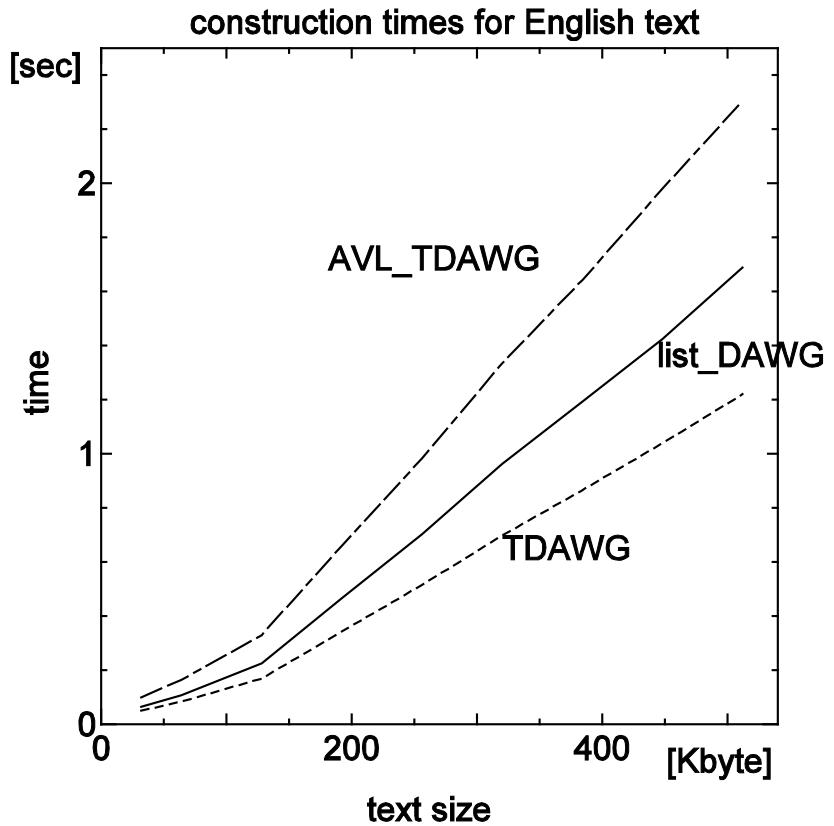
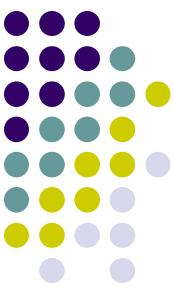
*For any string w of length n ,
 $\text{AVL_TDAWG}(w)$ can be constructed **on-line**
and in $O(n \log |\Sigma|)$ time and $O(n)$ space.*

Complexities

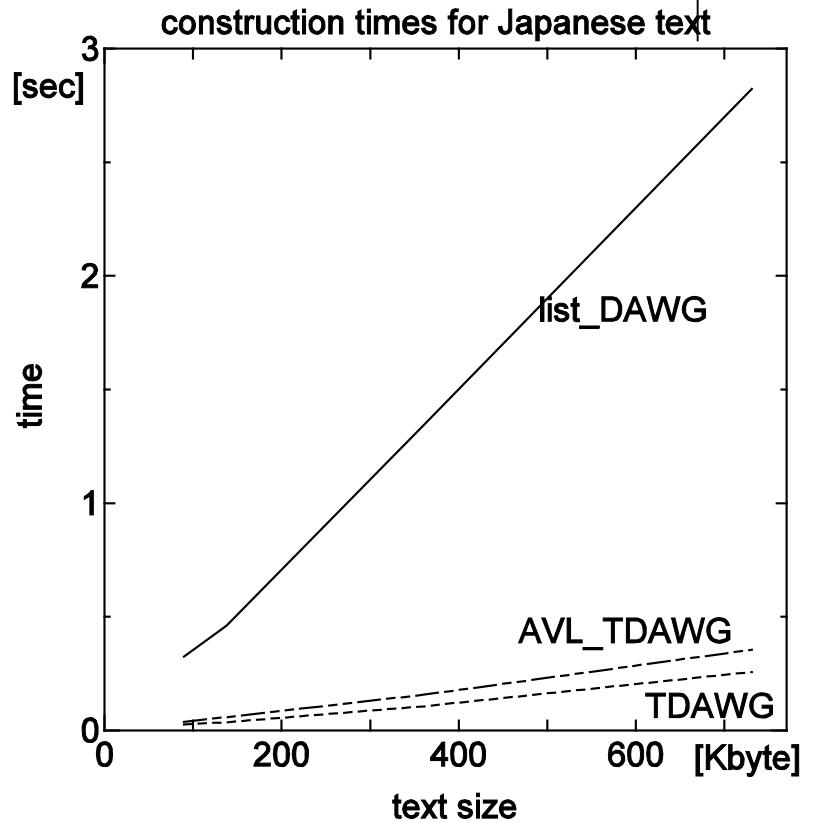


	space	search time	
		average case	worst case
<i>table</i>	$O(w \cdot \Sigma)$	$O(p)$	$O(p)$
<i>linked list</i>	$O(w)$	$O(p \cdot \Sigma)$	$O(p \cdot \Sigma)$
<i>TDAWG</i>	$O(w)$	$O(p \cdot \log \Sigma)$	$O(p \cdot \Sigma)$
<i>AVL-TDAWG</i>	$O(w)$	$O(p \cdot \log \Sigma)$	$O(p \cdot \log \Sigma)$

Experiment – Construction Times

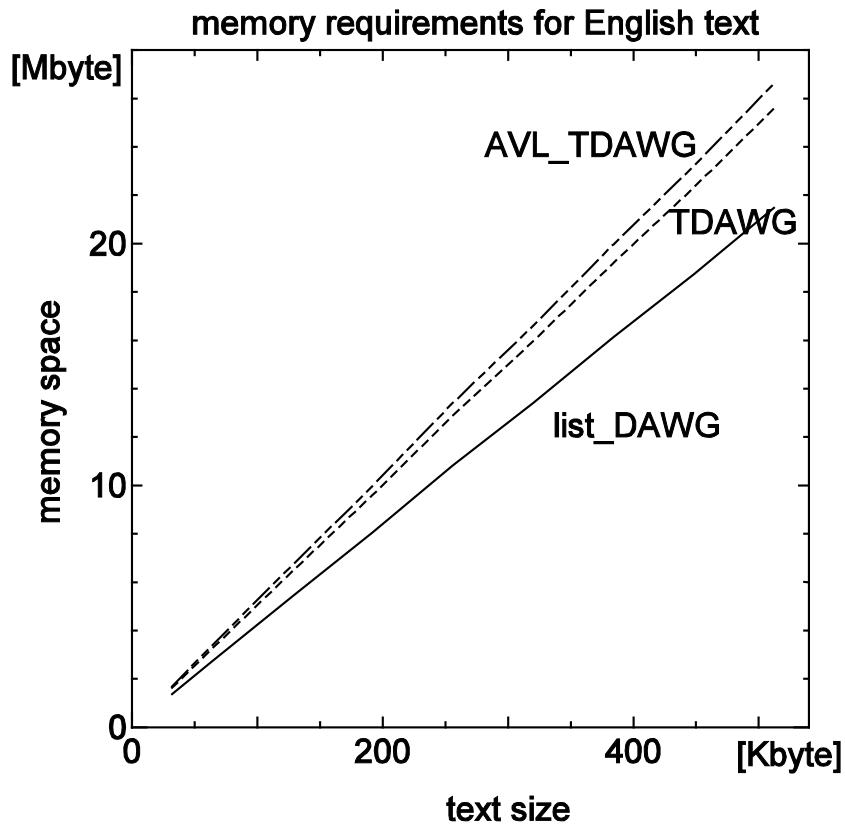
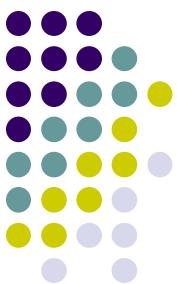


English

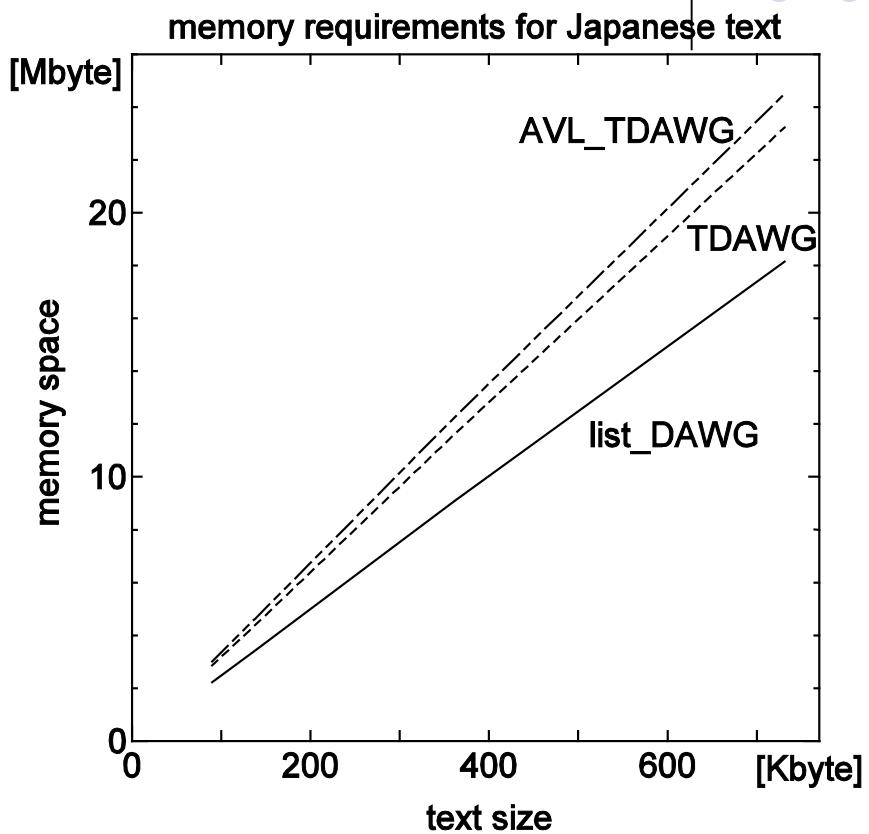


Japanese

Experiment – Memory Space

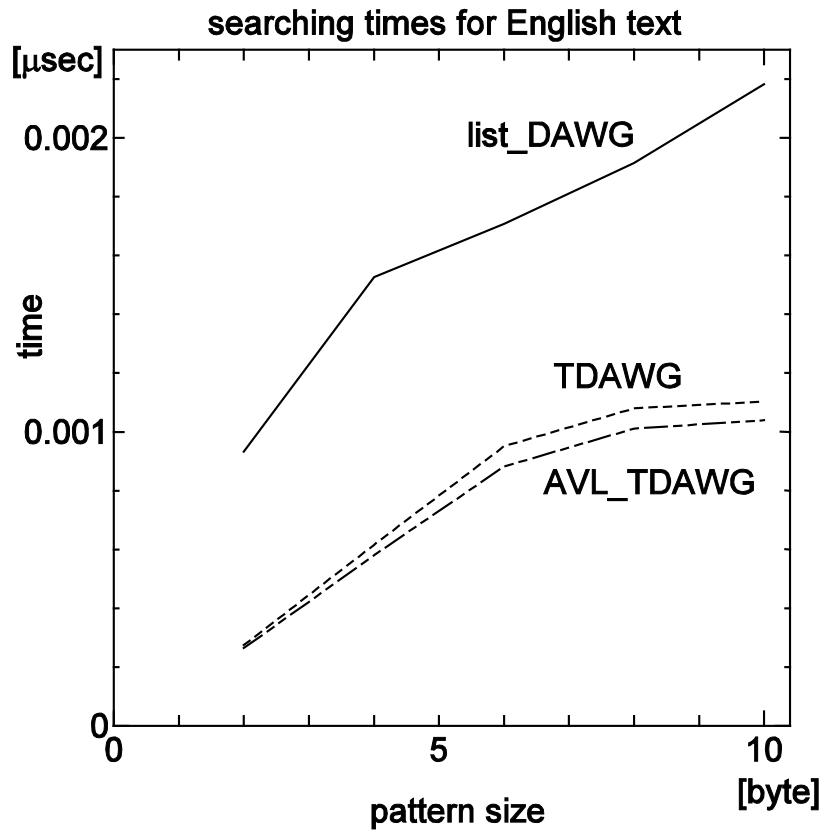
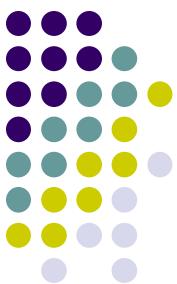


English

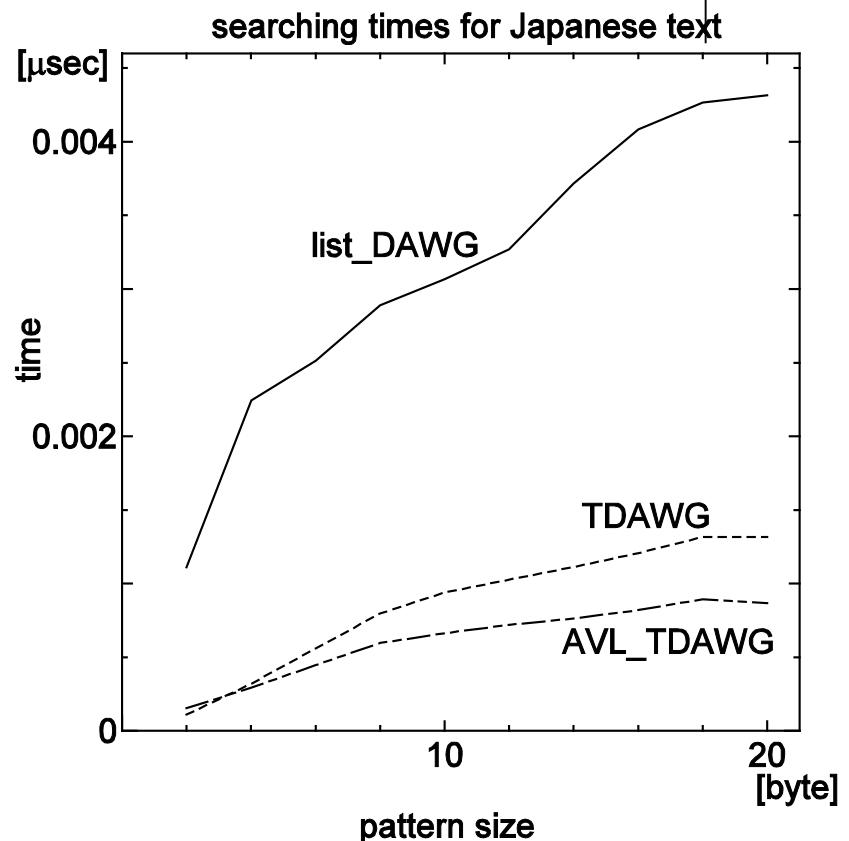


Japanese

Experiment – Search Times



English



Japanese



Citations

A. Blumer, J. Blumer, D. Haussler, A. Ehrenfeucht,
B. M. Chen, and J. Seiferas.

The smallest automaton recognizing the subwords of a text.
Theoretical Computer Science, 40:31-55, 1985.

M. Crochemore,

Transducers and repetitions.

Theoretical Computer Science, 45(1): 63-86, 1986

S. Miyamoto, S. Inenaga, M. Takeda, and A. Shinohara

Ternary Directed Acyclic Word Graphs.

8th Int. Conference on Implementation and Application of Automata
(CIAA 2003), Springer-Verlag LNCS 2759, pp. 121-130.
Also Accepted to Theoretical Computer Science