

Top-k String Similarity Search with Edit-Distance Constraint

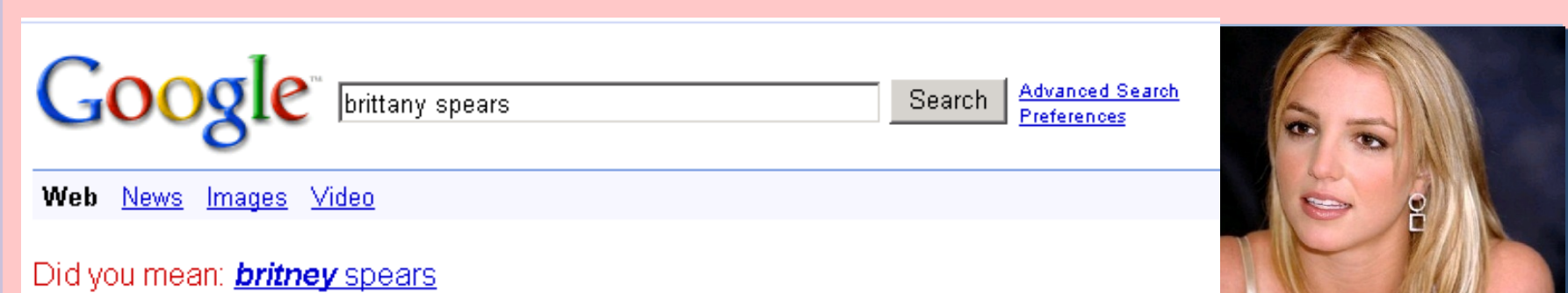
Dong Deng[§], Guoliang Li[§], Jianhua Feng[§], Wen-Syan Li[^]

[§]Department of Computer Science, Tsinghua University, Beijing, China

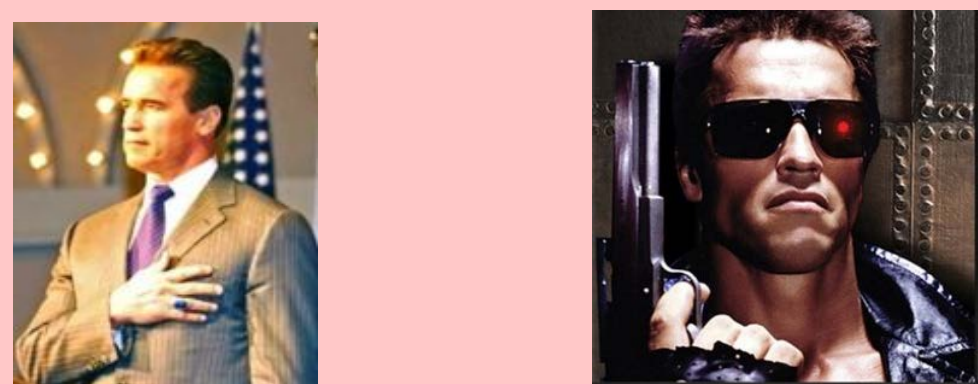
[^]SAP Labs, Shanghai, Beijing

Motivation

2008
Aravlis Zmils, Stephen P. Boyd, Dmitry M. Gorinevsky: Mixed state estimation for a linear Gaussian Markov model. IJCV 2008:3212-3226



The user doesn't know the exact spelling!



Find movies with a star "similar to" Schwarzenegger

Star	Title	Year	Genre
Keanu Reeves	The Matrix	1999	Sci-Fi
Samuel Jackson	Iron man	2008	Sci-Fi
Schwarzenegger	The Terminator	1984	Sci-Fi
Samuel Jackson	The man	2006	Crime

Top-k String Similarity Search

#1: Data in real world is dirty

Edit Distance: minimum # of single character transformations
e.g. $ED(\text{srajit}, \text{seraji}) = 2$

#2: Hard to define a threshold

#3: Many real applications

- Information retrieval
- Molecular biology
- Bioinformatics
- Data Quality, Data Cleaning

Problem Definition

Top-k String Similarity Search: Given a string set S and a query string q , top-k string similarity search returns a string set $R \subseteq S$ such that $|R|=k$ and for any string $r \in R$ and $s \in S - R$, $ED(r, q) \leq ED(s, q)$.

TABLE I

A STRING SET S AND A QUERY $q = \text{"srajit"}$

ID	s_1	s_2	s_3	s_4	s_5	s_6
String	sarit	seraji	suijt	suit	surajit	thrifty

the top-3 similar strings of *srajit*

Progressive Framework

(a) Traditional method

	0	1	2	3	4	5	6
0	0	1	2	3	4	5	6
1 s	1	0	1	2	3	4	5
2 e	2	1	1	2	3	4	5
3 r	3	2	1	2	3	4	5
4 a	4	3	2	1	2	3	4
5 j	5	4	3	2	1	2	3
6 i	6	5	4	3	2	1	2

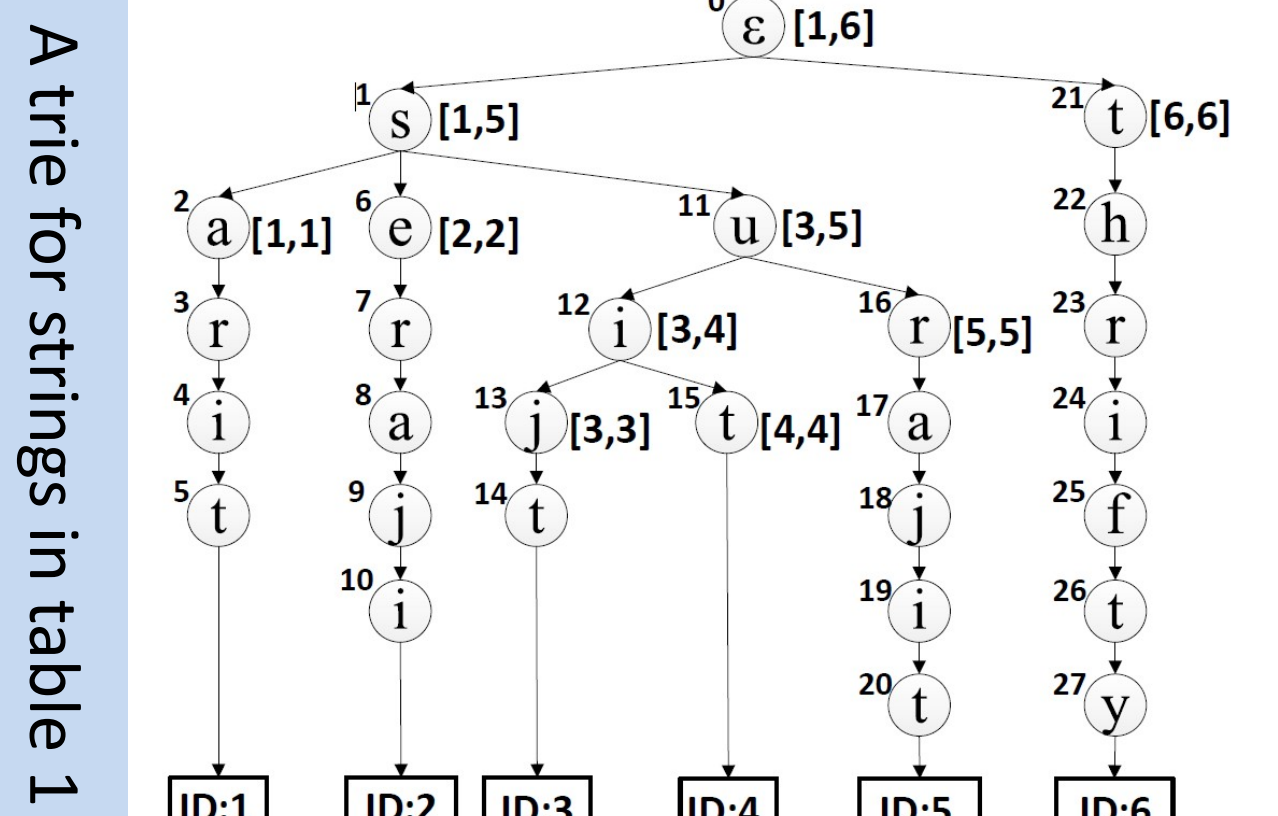
(b) Progressive method

	0	1	2	3	4	5	6
0	0	1	2				
1 s	1	0	1	2			
2 e	2	1	1	2			
3 r	3	2	1	2			
4 a	4	3	2	1	2		
5 j	5	4	3	2	1	2	
6 i	6	5	4	3	2	1	2

Traditional Method: Dynamic Programming

$D_{i,0} = i, D_{0,j} = j$, Insertion Deletion Match/Substitution

$D_{i,j} = \min\{D_{i-1,j} + 1, D_{i,j-1} + 1, D_{i-1,j-1} + 0/1\}$



Top-3 Query: *srajit*

(n_i, j) $\left\{ \begin{array}{l} i\text{-th node of trie;} \\ j\text{-th character of } Q \end{array} \right.$
Tx: node and char with $ED=x$

(a) $E_0 = \text{FINDMATCH}(-1, -1) = \{(0,0), (1,1)\}$
(b) Computing E_1 based on E_0

	Substitution	Insertion	Deletion	Substitution
EXTENSION	$\langle -1, 1 \rangle$	$\langle 1, 0 \rangle$	$\langle 0, 1 \rangle$	$\langle 2, 2 \rangle$

(c) Computing E_2 based on E_1

	Substitution	Insertion	Deletion	Substitution
EXTENSION	$\langle 2, 1 \rangle$	$\langle 1, 2 \rangle$	$\langle 3, 3 \rangle$	$\langle 3, 2 \rangle$
	$\langle 2, 0 \rangle$	$\langle 1, 1 \rangle$	$\langle 3, 2 \rangle$	$\langle 3, 1 \rangle$
	$\langle 1, 1 \rangle$	$\langle 0, 2 \rangle$	$\langle 2, 3 \rangle$	$\langle 2, 2 \rangle$

Progressive Method: Smallest Cell First
Two Operations: Find Match / Extend

(a) $T_0 = (n_0, 0) \cup \text{FINDMATCH}(0, 0) = \{(n_0, 0), (n_1, 1)\}$
(b) Computing T_1 based on T_0

	Substitution	Insertion	Deletion
EXTENSION	$\langle n_{1-1}, 1 \rangle$	$\langle n_{1-1}, 0 \rangle$	$\langle n_0, 1 \rangle$

(c) Computing T_2 based on T_1

	Substitution	Insertion	Deletion
EXTENSION	$\langle n_{21}, 1 \rangle$	$\langle n_{21}, 0 \rangle$	$\langle n_{20}, 1 \rangle$
	$\langle n_{22}, 2 \rangle$	$\langle n_{21}, 1 \rangle$	$\langle n_{21}, 0 \rangle$
	$\langle n_{23}, 2 \rangle$	$\langle n_{22}, 1 \rangle$	$\langle n_{22}, 0 \rangle$
	$\langle n_{21}, 2 \rangle$	$\langle n_{20}, 2 \rangle$	$\langle n_{21}, 1 \rangle$

Pivotal Entry-based Method

Definition(Pivotal Entry): An entry $\langle i, j \rangle$ in E_x is called a pivotal entry, if $D[i+1][j+1] > D[i][j]$.

S

	0	1	2	3	4	5	6
0	0						
1 s	0	1	2				
2 e	1	1	2				
3 r			2				
4 a				1			
5 j					1		
6 i						1	2

TABLE IV
PIVOTAL ENTRIES TO COMPUTE EDIT DISTANCE("srajit", "seraji")

(a) $E_0^p = \{(1, 1)\}$
(b) Computing E_1^p based on E_0^p

	Substitution	Insertion	Deletion
EXTENSION	$E_1^p[0] = \langle 2, 2 \rangle$	$E_1^p[1] = \langle 2, 1 \rangle$	$E_1^p[-1] = \langle 1, 2 \rangle$

(c) Computing E_2^p based on E_1^p

	Substitution	Insertion	Deletion
EXTENSION	$E_2^p[0] = \langle 3, 3 \rangle$	$E_2^p[1] = \langle 3, 2 \rangle$	$E_2^p[-1] = \langle 2, 3 \rangle$

Definition (Pivotal Triple): Given an entry $\langle n, j \rangle$, one of n 's children n_c , and a query q , triple $\langle n, j, n_c \rangle$ is called a pivotal triple, if $ED(n_c, q[1, j+1]) > ED(n, q[1, j])$.

(a) $T_0^p = \{(n_0, 0, n_{21}), (n_1, 1, n_2), (n_1, 1, n_6), (n_1, 1, n_{11})\}$
(b) Computing T_1^p based on T_0^p

	Substitution	Insertion	Deletion
EXTENSION	$T_1^p[0] = \langle n_{21}, 1, n_{22} \rangle$	$T_1^p[1] = \langle n_{21}, 0, n_{22} \rangle$	$T_1^p[-1] = \langle n_0, 1, n_{21} \rangle$

(c) Computing T_2^p based on T_1^p

	Substitution	Insertion	Deletion
EXTENSION	$T_2^p[0] = \langle n_{22}, 2, n_{23} \rangle$	$T_2^p[1] = \langle n_{22}, 1, n_{23} \rangle$	$T_2^p[-1] = \langle n_{21}, 2, n_{22} \rangle$

Range-based Method

Definition 4 (Pivotal Quadruple): A quadruple $\langle [l, u], d, j \rangle$ is a pivotal quadruple, if it satisfies (1) $\langle l, u \rangle$ is a sub-range of a d -th level node's range; (2) for any string s with ID in $[l, u]$, $ED(s[1, d+1], q[1, j+1]) > ED(s[1, d], q[1, j])$; (3) strings with ID $l-1$ or $u+1$ do not satisfy conditions (1) or (2).

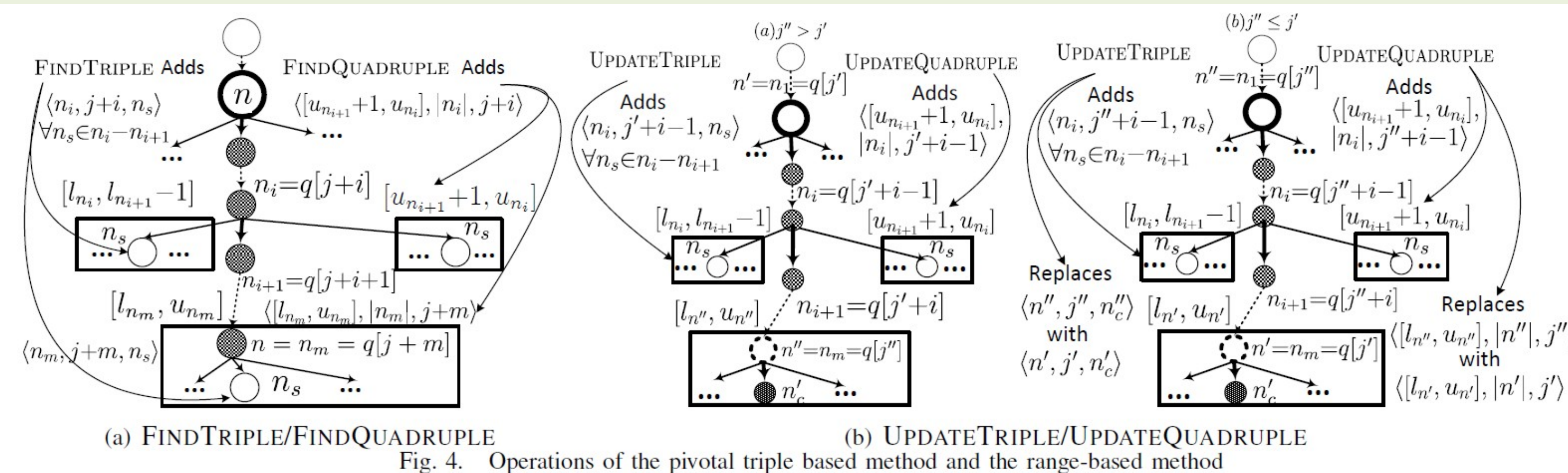


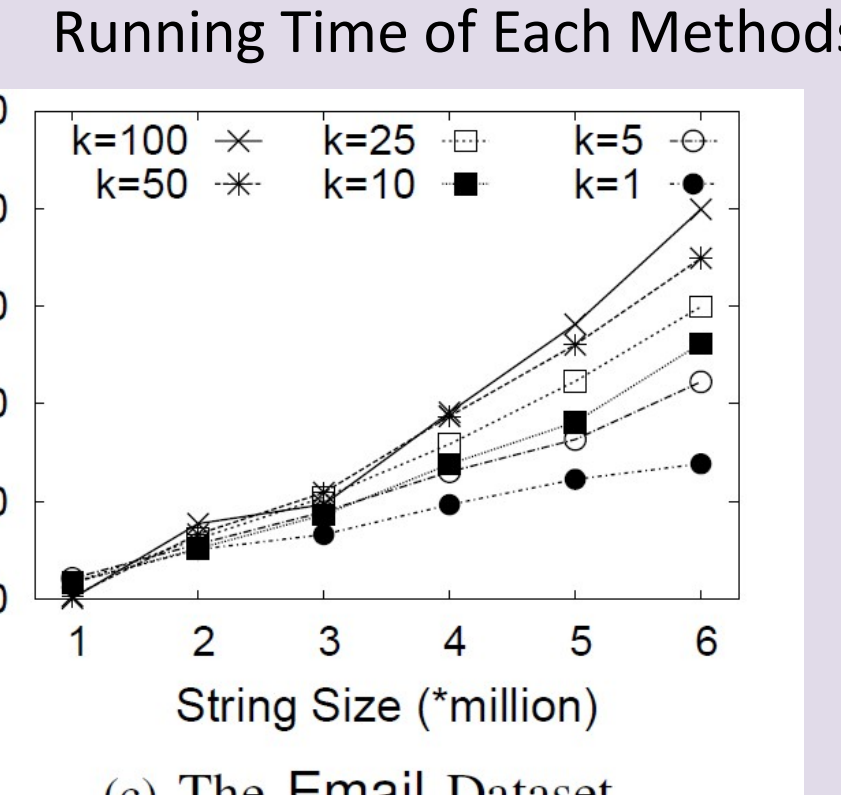
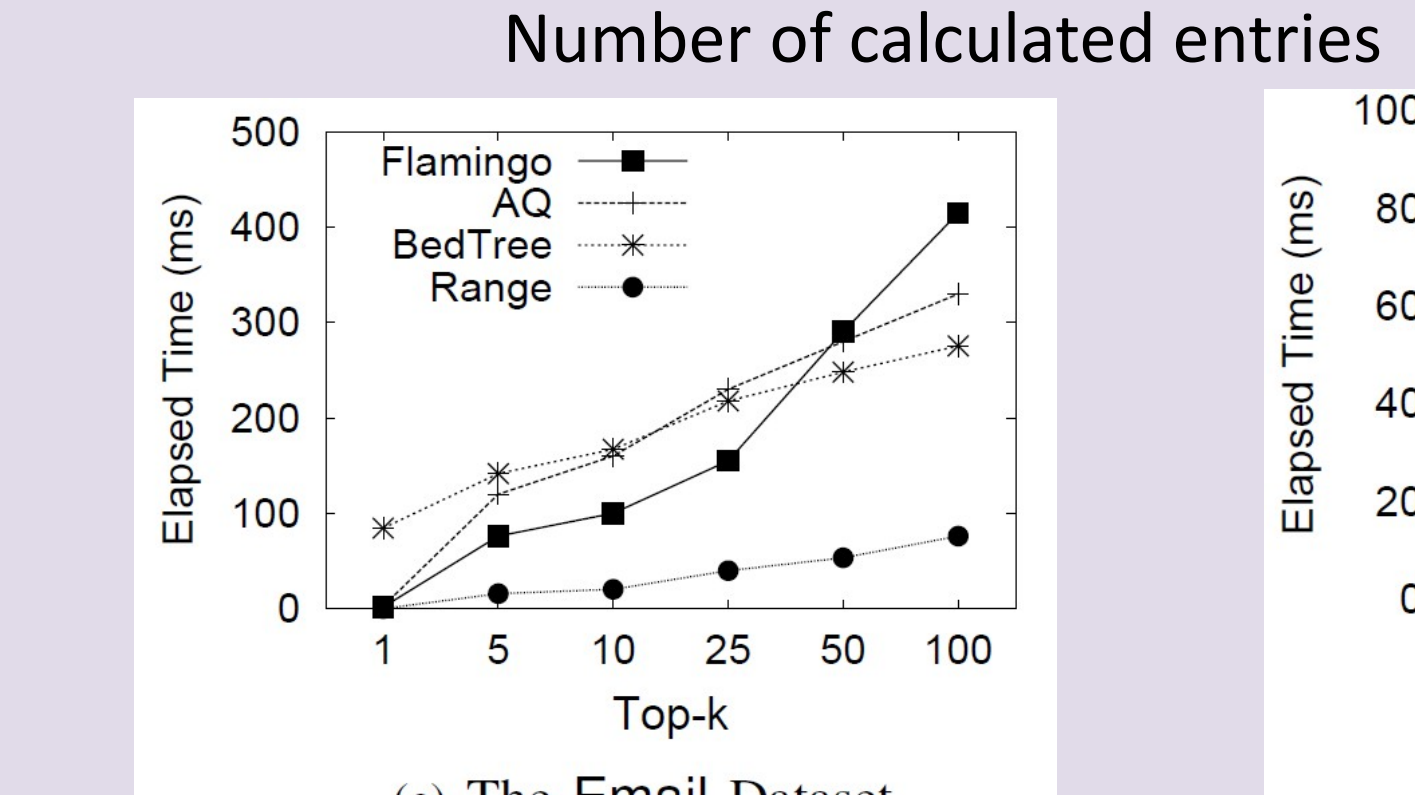
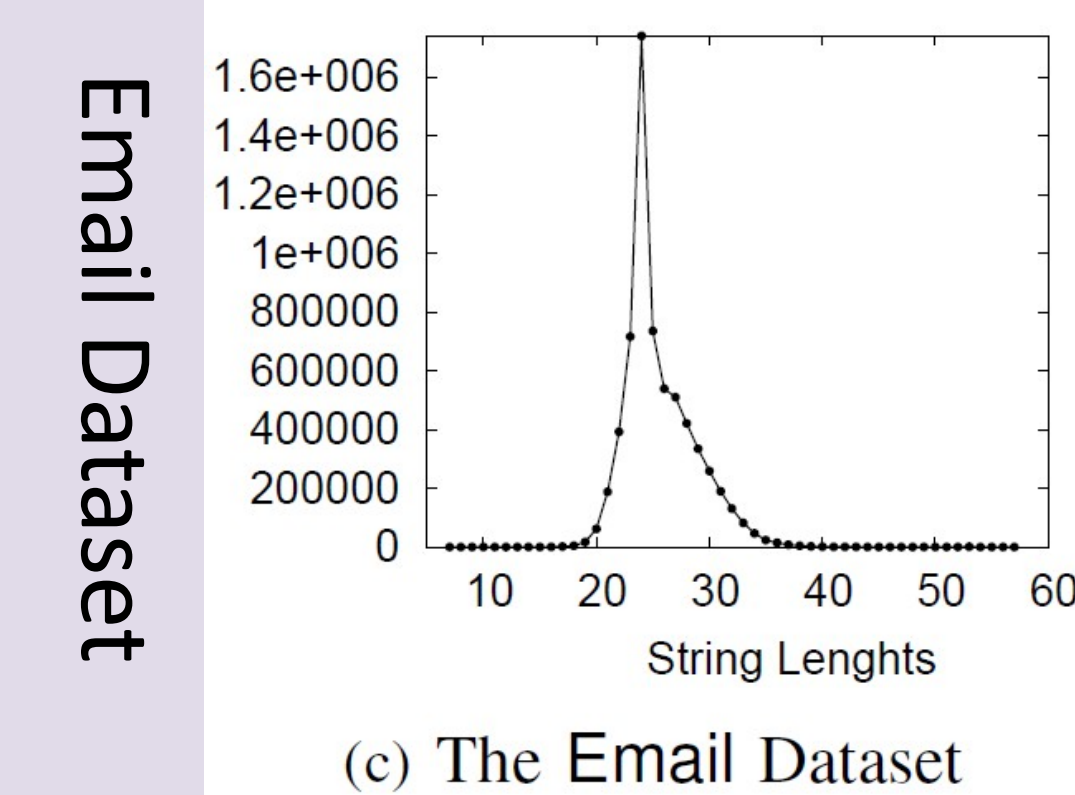
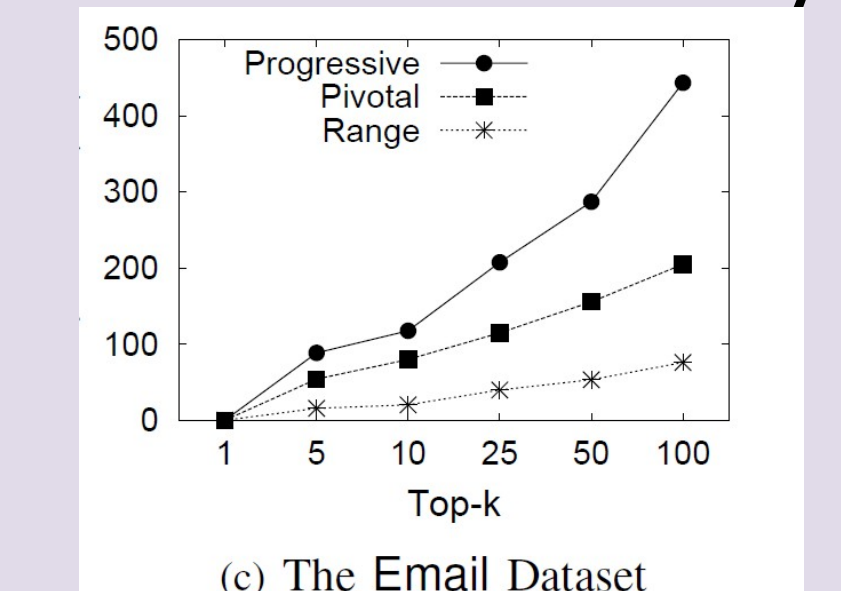
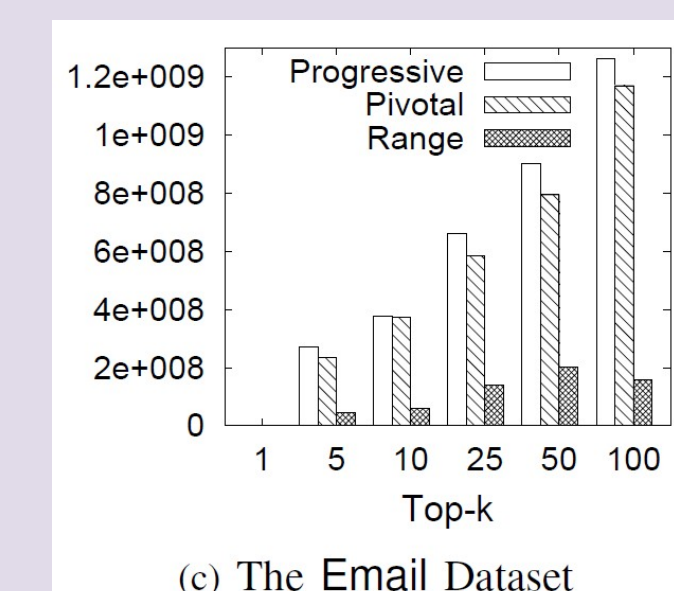
Fig. 4. Operations of the pivotal triple based method and the range-based method

Experiments

implemented in C++, Ubuntu: Intel Xero X5670 2.5GHz CPU and 4 GB memory

TABLE VII
DATASETS

Datasets	Cardinality	Avg Len	Max Len	Min Len
Word	146,033	16.01	35	1
Author	10.27 million	22.02	383	8
Email	6.4 million	26.58	57	7



Length Distribution State-of-the-art methods Scalability
<http://dbgroup.cs.tsinghua.edu.cn/dd/projects/topksearch>