

# A Unified Probabilistic Framework for Name Disambiguation in Digital Library

Jie Tang, A.C.M. Fong, Bo Wang, and Jing Zhang

**Abstract**—Despite years of research, the name ambiguity problem remains largely unresolved. Outstanding issues include how to capture all information for name disambiguation in a unified approach, and how to determine the number of people  $K$  in the disambiguation process. In this paper, we formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people  $K$ . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number  $K$  automatically found by our method is close to the actual number.

**Index Terms**—Digital libraries, information search and retrieval, database applications, heterogeneous databases.

## 1 INTRODUCTION

DIFFERENT people have the same name, and the same person has different names. In a digital library, the name ambiguity problem is a common problem. For example, in the DBLP database, there are 114 people with the name "Jie Tang". The number of people with the same name is a significant problem for name disambiguation. In this paper, we formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people  $K$ . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number  $K$  automatically found by our method is close to the actual number.

### 1.1 Motivation

We believe that the name ambiguity problem is a common problem in a digital library. For example, in the DBLP database, there are 114 people with the name "Jie Tang". The number of people with the same name is a significant problem for name disambiguation. In this paper, we formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people  $K$ . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number  $K$  automatically found by our method is close to the actual number.

- J. Tang and J. Zhang are with the Department of Computer Science and Technology, Tsinghua University, Rm 1-308, FIT Building, Beijing 100084, China. E-mail: jietang@tsinghua.edu.cn, zhangjing0544@gmail.com.
- A.C.M. Fong is with the School of Computing and Mathematical Sciences, Auckland University of Technology, AUT Tower Level 1, 2-14 Wakefield Street, Auckland 1142, New Zealand. E-mail: afong@aut.ac.nz.
- B. Wang is with the Department of Computer Science, Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China. E-mail: bowang@nuaa.edu.cn.

Manuscript received 1 July 2008; revised 5 Apr. 2010; accepted 16 Nov. 2010; published online 27 Dec. 2010.

Recommended for acceptance by B.C. Ooi.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-2008-07-0335. Digital Object Identifier no. 10.1109/TKDE.2011.13.

different people have the same name, and the same person has different names. In a digital library, the name ambiguity problem is a common problem. For example, in the DBLP database, there are 114 people with the name "Jie Tang". The number of people with the same name is a significant problem for name disambiguation. In this paper, we formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people  $K$ . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number  $K$  automatically found by our method is close to the actual number.

### 1.2 Prior Work

The name ambiguity problem is a common problem in a digital library. For example, in the DBLP database, there are 114 people with the name "Jie Tang". The number of people with the same name is a significant problem for name disambiguation. In this paper, we formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people  $K$ . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number  $K$  automatically found by our method is close to the actual number.



TABLE 1  
Attributes of Each Publication  $p_i$

Attribute	Description
$p_i.title$	title of $p_i$
$p_i.pubvenue$	published conference/journal of $p_i$
$p_i.year$	published year of $p_i$
$p_i.abstract$	abstract of $p_i$
$p_i.authors$	authors name set of $p_i$ $\{a_i^{(0)}, a_i^{(1)}, \dots, a_i^{(n)}\}$
$p_i.references$	references of $p_i$

## 2 PROBLEM FORMALIZATION

## 2.1 Definitions

I , e d c , a f , e a a b e  
 eaç a e p; a , Tab e 1. S ç b ca da a  
 ca be e ac ed f ce ç a DBLP, L b a. a.c ,  
 A e e . , a d C e ee . . . ed .

**Definition 1 (Principle Author and Secondary Author).**

Each paper  $p_i$  has one or more authors  $A_{p_i} = \{a_i^{(0)}, a_i^{(1)}, \dots, a_i^{(u)}\}$ .

We describe the author name that we are going to disambiguate as the principle author  $a_i^{(0)}$  and the rest (if any) as secondary authors.

We define the effective decay constant,  $\lambda_{\text{eff}}$ , by

- C P bVe e e (r<sub>1</sub>) e e e a e b, ed a e a e e e. F e a e, f b a e a e b, ed a "NKDD," e c e a e a d e c e d C P bVe e e a be e e e a e. I e, e e a c e e a e a e a d f f e e e e a c f e d, d b a e a d f f e e e e.
- C A (r<sub>2</sub>) e e e a a e p<sub>1</sub> a d p<sub>2</sub> a e a e c d a a e a e a e, e.e., A' p<sub>1</sub> ∩ A' p<sub>2</sub> ≠ ∅, e e A' p<sub>1</sub> d e e e f a f a e p<sub>1</sub> e c d e c e a a<sub>i</sub><sup>(0)</sup>, e.e., A' p<sub>1</sub> = A<sub>p1</sub> \ a<sub>i</sub><sup>(0)</sup>. T c a a e a a e a c c a d b e e a e e.
- C a (r<sub>3</sub>) e e e e a e c a e a e. I e a a a c e e e. F e, e c a e a e c a f a a f : If a e p<sub>1</sub> c e a e p<sub>2</sub>, p<sub>3</sub>, ..., p<sub>n</sub>, e e e a b d e c e d a e e a a a c e d a e, add d e c e d a e e a be e e p<sub>1</sub> a d e c e d a e.
- C a (r<sub>4</sub>) d e e c a e d a e feedbac. F a c e, e e c a e c f a a e d b e d a b a e d e a e e d b e d f f e e e.
- τ-C A (r<sub>5</sub>) e e e τ-e e C A e a. We e a e a e e a e a. S e a e p<sub>i</sub>, a a, "Nda d M c e, a d "Nda d e Ma, a d p<sub>j</sub>, a a, "Nda d M c e, a d "Nfe a d M f d. We a e d a b a e "Nda d M c e. A d f "Nda d e Ma, a d "Nfe a d M f d a c a a e a e, e e a p<sub>i</sub> a d p<sub>j</sub>, a e a 2-C A e a.

TABLE 2  
Relationships between Papers

$R$	$W$	Relation Name	Description
$r_1$	$w_1$	CoPubVenue	$p_i.pubvenue = p_j.pubvenue$
$r_2$	$w_2$	CoAuthor	$\exists r, s>0, a_i^{(r)}=a_j^{(s)}$
$r_3$	$w_3$	Citation	$p_i.cites\ p_j$ or $p_i.cites_s\ p_j$
$r_4$	$w_4$	Constraint	feedback supplied by users
$r_5$	$w_5$	$\tau$ -CoAuthor	$\tau$ -extension co-authorship ( $\tau>1$ )

$T_a e c e a, e e a f, e a b$   
 $d e e e, e e a e, a e a \tau-C A, e a$   
 $. F e e e a e d a a e, e c a c c a$   
 $c a e, e e e a c d e d e a a$   
 $a e a d e a c e d e d e a c a, e a . F$   
 $a a e p_1 a d p_2, e c a b a, e c e d-$   
 $e A'_{p_1} a d A'_{p_2} b, e c a .$  If  $a d f$   
 $A'_{p_1} \cap A'_{p_2} \neq \emptyset, e a, e a e, a e a C A,$   
 $e a . F d e e a 2-e e C A,$   
 $e a, e c c c a, e A^2_{p_1} a d A^2_{p_2}$   
 $a c c d, e c a, e . S e c f c a, A^2_{p_1}, e$   
 $e f a, b e e d A'_{p_1}, a e, b f, e$   
 $a, A'_{p_1}, e., A^2_{p_1} = A'_{p_1} \cup \{NB(a)\}_{a \in A'_{p_1}}, e e NB(a)$   
 $, e e f e, b f d e a. T e, e a, e$   
 $a e p_1 a d p_2, a e a 2-C A, e a, f a d$   
 $f A^2_{p_1} \cap A^2_{p_2} \neq \emptyset. F d e e, e e a e$   
 $, a e a 3-e e C A, e a, e f, e$   
 $e e d A^2_{p_1} f d a a, e A^3_p f e a c a e a d f$   
 $, e e, a e a e e c, e a, e a e$   
 $, a e a 3-C A, e a . T e e, f e a c e f$   
 $e a, r_i d e e d b w_i. E a, f, e a e f$   
 $d f f e e, b e d e c b e d$  Sec 4.  
 $I, e a e d a b a b e, e a e a$   
 $e a b e c e e d e, e a b e a e d e, e b$   
 $, e e. T e e a e b e a e d, e$   
 $d a b a a .$  We  $d e c b e c$  f  
 $a e a$  *cluster atom*.

**Definition 2 (Cluster Atom).** A cluster atom is a cluster in which papers are closely connected (e.g., the similarity  $K(x_i, x_j) > \text{threshold}$ ). Papers with similarity less than the threshold will be assigned to disjoint cluster atoms.

F d c e a d b e e a e f a e  
d a b a . F e a e , e c a a e , e c e a  
a , e a a f , e d a b a a , . F  
f d , e c e a , e c a e a c a e d-ba e d  
c e a e e c a . I  
add , e d e f e , e c c e f *cluster centroid*. De e d  
f , e c e a a , e e a e c a  
e , d f d , e c e d f a c e , e d a a  
a e a e , e c e f , e c e e c e d , a  
c a c a e d a , e a e c e a f a d a a  
a e d , e c e .

## 2.2 Name Disambiguation

G e a e a e a, e de e b ca c a  
 , e a a e a a  $P = \{p_1, p_2, \dots, p_n\}$ . T e b ca  
 da a e a ca be de ed b e  
 c de a d ed e . We e a ada e e

f, e -ca ed f a e a, [13] e e e, e  
b ca da a. P b ca a d e a, a e a -  
f ed a d e ced a, ç eac de  
e e e a a e a deac ed e a e a, . A b e  
f a a e a e a a c ed, e c e d de a a  
fea e ec . F, e ec, e e d (afe  
d f e a d e ) , e a b e f a a e a  
fea e a d e, e be f, e cc e ce a, e  
a e . F a, e ca def e, e b ca f a e  
a, a f :

**Definition 3 (Publication Informative Graph).** Given a set of papers  $P = \{p_1, p_2, \dots, p_n\}$ , let  $r_k(p_i, p_j)$  be a relationship  $r_k$  between  $p_i$  and  $p_j$ . A publication informative graph is a graph  $G = (P, R, V_P, W_R)$ , where each  $v(p_i) \in V_P$  corresponds to the feature vector of paper  $p_i$  and  $w_k \in W_R$  denotes the weight of relationship  $r_k$ . Let  $r_k(p_i, p_j) = 1$  iff there is a relationship  $r_k$  between  $p_i$  and  $p_j$ ; otherwise,  $r_k(p_i, p_j) = 0$ .

S e, e e a e K e  $\{y_1, \dots, y_K\}$ , e a e  
a, a d a b a e, e n b ca, e ea  
e ea ç e  $y_i, i \in [1, K]$ . M e ec f ca, e a a f  
a e d a b a ca be def ed a :

1. F a, e d a b a be . T e f -  
a a eed c de b, ca a b e  
fea e a ca ed, eac a e a d ea -  
be ee a e .
  2. S e be a c ed a a c . Ba ed  
e f a a, e a c ed a a c  
a d e a eff ce a .
  3. Dee e be f e e K. G e a  
d a b a a ( a f a -  
) , de e e, eac a K.
- I a e f, e e a . F,  
ed a e cea, f a e, e e d a b -  
a be a fed f a e . Sec d, e  
a, de, e., Ma Ra d Fed [16], a e a  
a ed de ea a da a. H e e, e  
b ca f a e a, e a e, be  
a b a c e ced b d ffe e e f ea,  
I cea, e f fe ce ( a a e e  
e a ) ç a a, a b a c e. I  
add, e a, e be f e e K a a  
ç a e a .

### 3 OUR FRAMEWORK

#### 3.1 Basic Idea

We, a e b a c b e a f, e a e d a b a -  
be : 1) a e, a c e ed, a e  
e a e a be (be e a e a, ); a d 2) a e  
a e a, e d, a e, e a e a be, f  
e a e, a e, c a, a a, a  
e a e . A dea d a b a e, e  
a e b e e a b, c e a a d a e  
ea, . T a a be, beca e  
e c e e, d ca e ba a ce, e  
ece f f a .  
I a e, e ea fed f a e ba ed  
Ma Ra d Fed [16], [24]. M e acc a e, e

f a e b, c e -ba ed f a a d c e -  
ba ed f a a H dde Ma Ra d Fed  
(HMRf) de a fea e f c . T e c b  
de ee f, e e f f a a e f a ed a  
e, f, e fea e f c . T e a ce f d ffe e  
e f ea, a de ed a e, f  
c e d fea e f c . S e HMRf de  
c de b, e a, e e, f fea e f c  
a da a e d ffe e e . S ç a f a e  
a ffe add a ad a a e : f,  
e ed ea, e ed ea, a d e -  
e ed ea . I a e, e f c  
e ed ea f a e d a b a, b  
ea c a e e / e ed f a  
e de . Sec d, a a d de eec, e  
HMRf de . T e bec e f c, e HMRf de  
a e bab d b f, dde a a be e  
b e a, ç a c e f de eec a e .

#### 3.2 Hidden Markov Random Fields

A Ma Ra d Fed a c d a bab  
d b f a be ( dde a a be ) , a be, e  
Ma e [16]. Ma eca ca e f MRF ca be  
de e ed. A H dde Ma Ra d Fed a e be  
f, e fa f MRF a d c ce de ed f  
H dde Ma M de (HMM) [15]. A HMRf a  
c ed f, e e c e : a be a be e f  
a d a a be  $X = \{x_i\}_{i=1}^n$ , a, dde fed f a d  
a a be  $Y = \{y_i\}_{i=1}^n$ , a d e, b, d be ee eac  
a f a a be, e, dde fed.  
We f a e, e d a b a be a, a f  
ea a a e d ffe e c e . Le, e  
dde a a be Y be, e c e a be, e a e .  
E e, dde a a be  $y_i$  a e a a e f, e e  
 $\{1, \dots, K\}$ , ç a e, e de e f, e c e . T e  
be a a a be X c e d a e, e e e e  
a d a a be  $x_i$  e e a ed f a c d a  
bab d b  $P(x_i|y_i)$  de e ed b, e c e -  
d, dde a a be  $y_i$ . F, e, e a d a a be  
X a e a ed be e e a ed c d a de e de  
f, e, dde a a be Y, e.,

$$P(X|Y) = \prod_{x_i \in X} P(x_i|y_i). \quad (1)$$

F . 2, e a, ca c e f, e HMRf f, e  
e a e F . 1. We ee, a de e de ed e a e  
ded be ee, e, dde a a be c e d  
e ea, F . 1. T e a e f eac, dde  
a a be (e.,  $y_1 = 1$ ) de e, ea e e . We d  
de, e d e c ea, be ee e, b,  
b, e de ca a a e, e de e de c e a,  
e ea,  
A HMRf a eca ca e f MRF, e bab  
d b f, e, dde a a be be, e Ma  
e . T, e bab d b f, e a e f  
 $y_i$  f, e be a a a be  $x_i$  de e d, e  
c e a be f b e a, a, a e e a,  $x_i$   
[24]. B, e f da e a, e e f a d fed [16],  
e bab d b f, e a be c f a Y  
a, e f

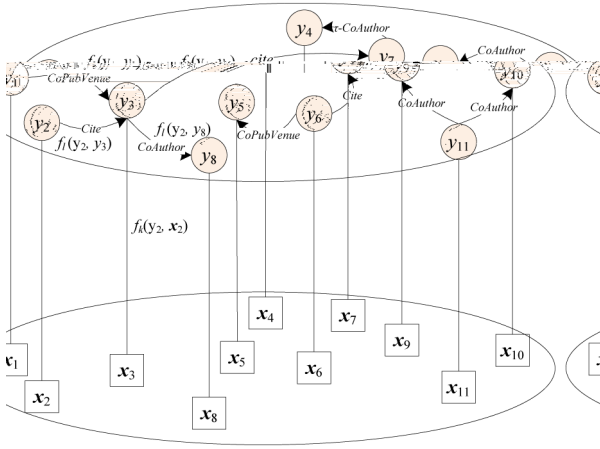


Fig. 2. Graphical representation of the HMRF model.  $f(y_i, y_j)$  and  $f(y_i, x_i)$  are edge feature and node feature, respectively, and will be described in the next section.

$$P(Y) = \frac{1}{Z_1} \exp \left( \sum_{(y_i, y_j) \in E, k} \lambda_k f_k(y_i, y_j) \right), \quad (2)$$

$$Z_1 = \sum_{y_i, y_j} \sum_{(y_i, y_j) \in E, k} \lambda_k f_k(y_i, y_j)$$

and  $b, f, e, e, c, e, b, ca, da, a, be, e, ea, ed, de, e, e, ca, Ga, a, d, b, e, a, e$

$$P(X|Y) = \frac{1}{Z_2} \exp \left( \sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i) \right), \quad (3)$$

$$Z_2 = \sum_{y_i} \sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i),$$

,  $e, e, f_k(y_i, y_j)$  a  $e, a, e, e, a, f, c$  (a  $ca, ed, e, fea, e, f, c$ )  $def, ed, ed, e, (y_i, y_j), a, d, E, e, e, e, a, ed, e, e, a, ; f_l(y_i, x_i)$  a  $e, a, f, c, def, ed, de, x_i; \lambda_k, a, d, \alpha_l, a, e, e, f, e, ed, e, fea, e, f, c, a, d, e, de, fea, e, f, c, e, ec, e; Z_1, a, d, Z_2, a, e, a, a, fac$ .

T  $fac, a, e, f, e, d, c, e, e, e, a, f, e, e, X, de, e, e, b, ca, e, P, a, d, e, x_i, de, e, e, ec, v(p_i), f, e, a, e, p_i$ .

### 3.3 Disambiguation Objective Function

We define a  $bec, e, f, c, a, e, Ma, a, P, e, c, f, a, f, e, HMRF, .e., b, a, P(Y|X), P(X)$  a  $a, a, e, a, c, a$ . T  $ee, f, e, acc, d, e, Ba, e, e, P(Y|X) \propto P(Y)P(X|Y)$ ,  $bec, e, f, c, ca, be, def, ed, a$

$$L_{\max} = \log(P(Y|X)) = \log(P(Y)P(X|Y)). \quad (4)$$

B  $b, (2), a, d, (3), (4), e, b, a$

$$L_{\max} = \log \left( \frac{1}{Z_1 Z_2} \exp \left( \sum_{(y_i, y_j) \in E, k} \lambda_k f_k(y_i, y_j) + \sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i) \right) \right). \quad (5)$$

$E, e, a, , e, ab, e, bec, e, f, c, , e, e, e, d, f, fea, e, f, c, , de, fea, e, f, c, f_l(y_i, x_i), a, d, ed, e, fea, e, f, c, f_k(y_i, y_j), e, e, e, e, a, b, e, f, a, a, ca, ed, eac, a, e, a, d, e, ea, f, a, be, ee, a, e, e, ec, e, . T, e, ed, e, fea, e, f, c, f_k(y_i, y_j), ed, ca, ace, e, e, ea, be, ee, a, e, . I, e, f, a, e, a, ea, ea, a, da, a, e, a, eac, e, e, e, e, a, e, a, e, be, a, ed, e, a, ec, e, . S, ec, f, ca, , e, ed, e, fea, e, f, c, ca, e, e, a, e, ea, ca, C, P, b, Ve, e, a, d, C, A, (a, Tab, e, 2), a, d, a, ea, e, f, a, . T, , e, def, e, e, ed, e, fea, e, f, c, a$

$$f_k(y_i, y_j) = K(x_i, x_j) \sum_{r_m \in R_{ij}} [w_m r_m(x_i, x_j)]. \quad (6)$$

He  $e, K(x_i, x_j)$  a  $a, a, f, c, be, ee, a, e, x_i, a, d, x_j; w_m, e, e, f, ea, r_m; R_{ij}, de, e, e, e, f, ea, be, ee, x_i, a, d, x_j; a, d, r(x_i, x_j), de, e, a, f, c, f, e, ea, be, ee, x_i, a, d, x_j. T, e, e, a, def, e, e, ea, f, c, r(x_i, x_j), def, e, b, a, a, e, a, de, c, bed, Def, 3. He, e, ef, e, c, de, a, def, c, c, b, e, e, e, f, a, .e., r_1(x_i, x_j) = \exp\{-|x_i.year - x_j.year|\}. T, def, de, ed, f, a, be, a, e, a, e, a, b, be, : e, C, A, a, d, C, P, b, Ve, e, ea, a, e, fe, e, de, e, de, .e., a, e, d, b, a, e, ea, f, c, ed, c, fe, e, ce, / a, e, e, a, ec, f, c, e, da, d, ca, a, e, d, c, ab, a, e, eac, e, a, ec, f, c, e, d.$

T, e, de, fea, e, f, c,  $f_l(y_i, x_i)$  a  $ca, e, e, a, b, e, f, a, a, ca, ed, a, e, x_i. T, e, ba, c, dea, e, e, f, e, a, e, a, a, e, e, a, e, a, c, e, e, e, e, a, e, a, e, be, a, ed, e, c, e. F, c, e, a, e, ea, e, def, e, e, de, fea, e, f, c, a$

$$f_l(y_i, x_i) = K(y_i, x_i) = K(\mu_{(i)}, x_i), \quad (7)$$

,  $e, e, \mu_{(i)}, e, c, e, ce, d, a, e, a, e, x_i, a, ed. N, a, K(x_i, \mu_{(i)})$   $e, e, e, e, a, be, ee, a, e, x_i, a, d, a, ed, c, e, ce, e, \mu_{(i)}. T, e, (6), a, d, (7), (5), e, b, a$

$$L_{\max} = \sum_{(x_i, x_j) \in E, k} \lambda_k K(x_i, x_j) r_k(x_i, x_j) + \sum_{x_i \in X, l} \alpha_l K(x_i, \mu_{(i)}) - \log Z, \quad (8)$$

,  $e, e, Z = Z_1 Z_2. W, a, f, e, ea, e, c, b, e, e, e, f, ed, e, fea, e, f, c, \lambda_k, a, d, e, e, f, e, ea, w_m, a, d, ea, \lambda, f, c.$

### 3.4 Criteria for Model Selection

We  $e, Ba, e, a, I, f, a, C, e, (BIC), a, ec, e, e, a, e, e, be, f, e, e, K. We, def, ea, bec, e, f, c, f, e, d, a, b, a, a, . O, a, e, a, a, a, e, e, a, a, e, e, ca, bec, ef, c, e, e, K, a, d, f, da, be, K, a, a, e, e, ba, bec, ef, c.$

Secfca, ef c de  $K=1$ , a, e e  
e e, e e a e a. Te, e e a  
ea e e de e e, e e, e a e c e  
d be bc e. Ne, f eac  
bc e, e a a e, e ea e e de e e  
e, e. Te ea e ea e c d-  
a fed (e., bc e ca be). I, e  
ce, e ca  $M_h$ , e de c e d, e  
e e be h. We, e ef e, a e a  
fa fa e a e de  $M_h$ , e e h a e f 1  
 $n$ , c e.  
N, a ç e, e be de f  $M_h$ .  
Ma ea e e ca be ed f de e ec,  
ç a S, e e C eff ce [23], M De c  
Le, (MDL) [34], A a e I f a C e (AIC)  
[1], a d e bab e a [22]. We ç e  
BIC a, e c e, beca e BIC c e f da e -  
a a, e c e a ç a MDL a d, a a  
e e a, a, e, e c e a ç a AIC,  
ç de abe be. Ba ed e e e  
c de a, e e a a a f, e BIC ea e e  
[22] a, e c e

$$BIC^v(M_h) = \log(P(M_h|P)) - \frac{|\lambda|}{2} \cdot \log(n), \quad (9)$$

e e  $P(M_h|P)$ , e e bab f de  $M_h$   
e, e be a  $P \cdot |\lambda|$ , e be f a a e e  
 $M_h$  (ç ca be def ed d ffe e a, e., e  
be f e a a e e, e de  $M_h$ , e  
f, e bab e f  $P(Y)$ ).  $n$ , e a e be. Te  
ec d a a e a de c e.  
I e e ce, a BIC c e a a e, a -  
a e, e de  $M_h$  f, e e da e e. We e,  
c e f, e de e ec beca e ca be ea  
e e ded d ffe e a. F e a e, c e -  
a c e a e  $K$ - ea [27] X-  
ea [33] e a d, e da a a de e de a d, e  
e bab  $P(M_h|P)$  ca be fed  
 $P(P|M_h)$  acc d, e Ba e a e  $P(M_h|P) \propto$   
 $P(P|M_h)P(W_h)$  e a e  $P(M_h)$  a d, e  
H e e, e e d a e ad a a e f de e de c e  
be ee, e c e e. Te, e  $P(M_h)$  a  
f a a e. O def (2) c de  
e de e de c e a Ma f e d.

## 4 PARAMETER ESTIMATION

### 4.1 Algorithm

Te a a e e e a be de e e, e  
a e f, e a a e e  $\Theta = \{\lambda_1, \lambda_2, \dots; \alpha_1, \alpha_2, \dots\}$  a d  
de e ea e f a a e. Me acc ae, e  
e, e - e, d bec ef c (8),  
e ec ac d a de  $P(Y|X, \Theta)$ .

A a, e e, e ea a (cf. A 1)  
f a a e e e a a c f ea e  
e: Assignment f a e, a d Update f a a e e  $\Theta$ .  
Te ba c dea, a e f a d ç e a  
a a e e e  $\Theta$  a d e ec a ce d f eac c e.  
Ne, ea eac a e c e c e a d, e  
ca c a e, e ce d f eac a e-c e ba ed, e

a e. Af e, a, e da e, e e, f eac  
fea ef c b a, e bec ef c.

### Algorithm 1: Parameter estimation

Input:  $P=\{p_1, p_2, \dots, p_n\}$

Output: model parameters  $\Theta$  and  $Y=\{y_1, y_2, \dots, y_n\}$ , where  $y_i \in [1, K]$

#### 1. Initialization

1.1 randomly initialize parameters  $\Theta$ ;

1.2 for each paper  $x_i$ , choose an initial value  $y_i$  with  $y_i \in [1, K]$

1.3 calculate each paper cluster centroid  $\mu_{(i)}$ ;

( $y_i, x_i$ ) 1.4 for each paper  $x_i$  and each relationship  $(x_i, x_j)$ , calculate  $f$   
and  $f_k(y_i, y_j)$ .

#### 2. Assignment

2.1 assign each paper to its closest cluster centroid;

#### 3. Update

3.1 update of each cluster centroid;

3.2 update of the weight for each feature function.

F a a, e a d a, e a e f eac  
a a e e ( $\lambda$  a d  $\alpha$ ). F a a f, e c e  
ce d, e f e a a, c e e, d  
de f, e c e a. Ba ca, a e, a  
e, a a, e, d be a ed d c e  
a. We ed a a e, e de c bed fa,  
b a a ç, e a e, a, a, e, e  
a, e c e ce d u. I, a, e e  
 $\gamma$  c e a. If  $\gamma$  e a, e be f e e  $K$ ,  
e, e e  $\gamma$  a e ed a a a e. If  
 $\gamma < K$ , e a d ç e a, e ( $K-\gamma$ ) a e a, e  
c e ce d. If  $\gamma > K$ , e e ea e c e  
a, e e a e  $K$  ef. We  
d ce de a, e e a a e e  
e a a.

**Assignments.** I Assignments, eac a e  $x_i$  a ed  
 $\mu_{(h)}$  a e  $\log P(y_i|x_i)$

$$\begin{aligned} \log P(y_i|x_i) &\propto L_{x_i}(\mu_{(h)}, x_i) \\ &= \sum_k \lambda_k K(x_i, x_j) r_k(x_i, x_j) \\ &\quad + \sum_l \alpha_l K(x_i, \mu_{(h)}) - \log Z, \end{aligned} \quad (10)$$

e e  $Z$  de ade a a a fac  $x_i$  a d  
ca be e ed a e ca e ab, e ea e c e

f c , e. , d e e a d a d e e. H e e ,  
 , e a e e e , e f c e a .  
 N , , e a c a c a e a a a e c e (10).  
 T e f e (10) a e a a c b a f  
 , e a f c  $K(x_i, \mu_{(h)})$  a d , e e a a -  
 a f c  $K(x_i, x_j)$ , , c a b e c a c a e d. H e e ,  
 a c a b e b a a e a c f , e a  
 f c , .e., (Z), b e c a e , e a a a d  
 a e a c e , , e a , (Z = Z<sub>1</sub>Z<sub>2</sub>). A f e a -  
 , , a e b e e e d f a a e f e e c e, e. ,  
 b e f a a [30] a d c a e d e e c e (CD)  
 [19]. W e e a e , a a e , e a  
 f c a c a e d e e c e d a b a  
 b e c e f c .  
 B a e d J e e ' e a [21], e c a b a a e  
 b d f , e e a e - e , d (L) , a K -  
 b a c -L e b e (KL) d e e c e

$$\begin{aligned} L^{KL} &= KL(q||P) \\ &= \sum_{y_i} q(y_i|x_i) \log(q(y_i|x_i)) - \sum_{y_i} q(y_i|x_i) \log(P(y_i|x_i)) \\ &= -H(q) - \langle \log(P(y_i|x_i)) \rangle_{q(y_i)}, \end{aligned} \quad (12)$$

, e e  $q(y_i|x_i)$  a a a f , e d b  
 $P(y_i|x_i) \cdot \langle \cdot \rangle_q$  , e e e c a d e , e d b q.  
 M a , e - e , d f , e d a a (5) e a -  
 e , e KL d e e c e (12) b e e e , e d a a  
 d b  $q^0$  a d , e e b d b e , e  
 b e a a b e ,  $q^\infty$  , e e , e f e c a b e c a c a e d  
 b , e b e a , e c e a e d a b e a d  
 , e e c d e , e b a b , e e e e d e  
 d b , a b e a b e . A a , , e b  
 d f f c , e a b e e a , d e e , e e c d  
 e . A M a c a M e C a (MCMC) e , d c a b e  
 e d e a e , e a a d b  $q^\infty(y_i|x_i)$   
 , e a f MCMC b e e e d a  $q^0(y_i|x_i)$ . T  
 a e , e c e e f f c e , e c a e , e c a e  
 d e e c e a , [19], , c a a e , e d -  
 b b e e a G b b a e ( e e ).  
 T , , e b e c e f c b e c e

$$\begin{aligned} L^{KL} &= KL(q^0||P) \approx KL(q^0||P) - KL(q^l||P) \\ &= \langle \log(P(y_i|x_i)) \rangle_{q^0(y_i)} - \langle \log(q^l(y_i|x_i)) \rangle_{q^l(y_i)}. \end{aligned} \quad (13)$$

I c a e d e e c e e a , e a d f  
 $KL(q^0||q^\infty)$ , e e , e d f f e e c e b e e e  $KL(q^0||q^l)$   
 a d  $KL(q^l||q^\infty)$ , , e e  $q^l$  , e d b e , e N - e ,  
 e c c f , e d a a e c (.e., b e a ) , a  
 a e e e a e d a f e l - e G b b a . A d c a e d  
 [19], e e l c a b e e a 1 c a e . (T a ,  
 e c a c d e e G b b a e a  
 e , e  $KL(q^0||q^l)$ ). T e c e d e f e c c  
 , e d a a e c (.e.,  $q^l$ ) f , e d b  $q^0$  d e c b e d  
 A , 2.

### Algorithm 2: One-step sampling

Input: current observation  $x^0$  and labels  $y^0$

Output: sampling results of  $y^1$  and  $x^1$

- 1: Draw an observation  $x_i$  from the distribution of  $q^0(x_i)$  ( $q(x)$  can be obtained by summing over all possible labels);
- 2: Compute  $P(y_i|x)$ , the posterior probability distribution over the label variable given the observation  $x_i$ ;
- 3: Compute  $P(y_i|y_{-i})$ , the probability distribution over the label variable given labels of its neighboring observations;
- 4: Draw a new label  $y_i^1$  for each observation from the probability distribution  $P(y_i|x)P(y_i|y_{-i})$ ;
- 5: Given the chosen label, compute the conditional distribution of  $P(x_i|y_i)$ .

F a , b a e d , e e c c e d d a a e c , e c a  
 c a c a e (13). T e c a c a e e e e  
 d e a d . T a e e e f f c e , e c a e , e  
 d e e c e a f e d a , [44] e a c e , e  
 a c e d e .  
 A f e , e , d e (10), e c a c e e , e  
 f , e , e b e c e f c . F a , a e e d  
 a , e d e e a d a e , e a e f  
 e a c a e . A a e f a a e e f e d , e  
 e e , e , e a e f e d. T e c e e e a e d  
 a e c a e a e b e e e  
 c c e e e a .

**Update.** I U d a e, e a c c e c e d f d a e d  
 b , e a , e c e a f , e a e c a e d

$$\mu(h) = \frac{\sum_{i:y_i=h} x_i}{\|\sum_{i:y_i=h} x_i\|_A}. \quad (14)$$

T e , b d f f e e a , e b e c e f c  
 e e c e a c a a e e  $\lambda_k$ , e , a e

$$\frac{\partial L}{\partial \lambda_k} = - \sum_{(x_i, x_j) \in E} K(x_i, x_j) r(x_i, x_j) - \frac{\partial \log Z}{\partial \lambda_k}. \quad (15)$$

W e e e , a , e e c d e a c a b e , b e c a e  
 c a c a f Z e e d a b e f  
 a e f e a c a e . A a , e a f , e KL  
 d e e c e b e c e f c (13) a d e , e CD a  
 c a c a e , e d e a e f  $L^{KL}$  , e e c  $\lambda_k$

$$\begin{aligned} \frac{\partial L^{KL}}{\partial \lambda_k} &= \left\langle \frac{\partial \log(P(y_i|x_i))}{\partial \lambda_k} \right\rangle_{q^0(y_i)} - \left\langle \frac{\partial \log(q(y_i|x_i))}{\partial \lambda_k} \right\rangle_{q^l(y_i)} \\ &= - \sum_{(x_i, x_j) \in E} K(x_i, x_j) r(x_i, x_j) - \left\langle \frac{\partial \log(q(y_i|x_i))}{\partial \lambda_k} \right\rangle_{q^l(y_i)}. \end{aligned} \quad (16)$$

T e f e , e a c b a f , e  
 a f c a d , e e c d e c a b e c a c a e d  
 a f e , e l - e a (A , 2).  
 F a , e a c a a e e d a e d b

$$\lambda_k^{new} = \lambda_k^{old} + \Delta \frac{\partial L}{\partial \lambda_k}, \quad (17)$$

, e e  $\Delta$  , e e a a e . W e d , e a e f  $\alpha$ .

## 4.2 Estimation of $K$

Our algorithm estimates  $K$  (see Algorithm 3) by iteratively finding the best two-sub-clusters model  $M_2$  for each cluster  $C$  in  $C^{(i)}$ . We calculate  $BIC(M_2)$  and  $BIC(M_1)$  for each cluster  $C$ . If  $BIC(M_2) > BIC(M_1)$ , we split cluster  $C$  into two sub-clusters  $C^{(i+1)} = \{C_1, C_2\}$ . We calculate  $BIC$  for the obtained new model, and repeat the process until no more splits are found.

### Algorithm 3. Estimation of $K$

Input:  $P = \{p_1, p_2, \dots, p_n\}$

Output:  $K, Y = \{y_1, y_2, \dots, y_n\}$ , where  $y_i \in [1, K]$

```

1:  $i=0, K=1$ , that is to view  $P$  as one cluster:  $C^{(0)} = \{C_1\}$ ;
2: do {
3:   foreach cluster  $C$  in  $C^{(i)}$  {
4:     find a best two sub-clusters model  $M_2$  for  $C$ ;
5:     if  $(BIC(M_2) > BIC(M_1))$ 
6:       split cluster  $C$  into two sub-clusters  $C^{(i+1)} = \{C_1, C_2\}$ ;
7:       calculate  $BIC$  score for the obtained new model;
8:   } while (existing split);
9: } while (no more splits found);

```

Our algorithm estimates  $K$  (see Algorithm 3) by iteratively finding the best two-sub-clusters model  $M_2$  for each cluster  $C$  in  $C^{(i)}$ . We calculate  $BIC(M_2)$  and  $BIC(M_1)$  for each cluster  $C$ . If  $BIC(M_2) > BIC(M_1)$ , we split cluster  $C$  into two sub-clusters  $C^{(i+1)} = \{C_1, C_2\}$ . We calculate  $BIC$  for the obtained new model, and repeat the process until no more splits are found.

$$\sum_{i=1}^K (P(y_i) + \mu_{(i)}) + \sum_{\lambda \in \Theta} \lambda. \quad (18)$$

## 5 EXPERIMENTAL RESULTS

### 5.1 Experimental Setting

**Data Sets.** We used two datasets: **ACM** [40] and **DBLP**. The ACM dataset contains 32 authors, 2,074 papers, and 1,000 citations. The DBLP dataset contains 40 authors, 4,000 papers, and 10,000 citations. We used the **ACM** dataset for the **ACM** experiments and the **DBLP** dataset for the **DBLP** experiments. We used the **ACM** dataset for the **ACM** experiments and the **DBLP** dataset for the **DBLP** experiments. We used the **ACM** dataset for the **ACM** experiments and the **DBLP** dataset for the **DBLP** experiments.

TABLE 3  
Data Sets

Abbr. Name	#Publications	#Actual Person	Abbr. Name	#Publications	#Actual Person
Cheng Chang	12	3	Gang Wu	40	16
Wen Gao	286	4	Jing Zhang	54	25
Yi Li	42	21	Kuo Zhang	6	2
Jie Tang	21	2	Hui Fang	15	3
Rakesh Kumar	61	5	Michael Wagner	44	12
Bing Liu	130	11	Jim Smith	33	5
Ajay Gupta	27	4	Wei Wang	306	90
Dimitry Pavlov	16	2	David Jensen	43	3
Charles Smith	7	4	David Brown	53	7
David C. Wilson	52	5	George Miller	17	2
James H. Anderson	112	2	James Johnson	17	3
John Miller	74	2	Joseph Miller	10	2
Paul Jones	13	3	Richard Taylor	93	10
Robert Fisher	105	4	Robert Moore	92	3
Robert Williams	8	2	William Cohen	110	2

Our algorithm estimates  $K$  (see Algorithm 3) by iteratively finding the best two-sub-clusters model  $M_2$  for each cluster  $C$  in  $C^{(i)}$ . We calculate  $BIC(M_2)$  and  $BIC(M_1)$  for each cluster  $C$ . If  $BIC(M_2) > BIC(M_1)$ , we split cluster  $C$  into two sub-clusters  $C^{(i+1)} = \{C_1, C_2\}$ . We calculate  $BIC$  for the obtained new model, and repeat the process until no more splits are found.

Our algorithm estimates  $K$  (see Algorithm 3) by iteratively finding the best two-sub-clusters model  $M_2$  for each cluster  $C$  in  $C^{(i)}$ . We calculate  $BIC(M_2)$  and  $BIC(M_1)$  for each cluster  $C$ . If  $BIC(M_2) > BIC(M_1)$ , we split cluster  $C$  into two sub-clusters  $C^{(i+1)} = \{C_1, C_2\}$ . We calculate  $BIC$  for the obtained new model, and repeat the process until no more splits are found.

**Experimental Design.** We used the **ACM** dataset for the **ACM** experiments and the **DBLP** dataset for the **DBLP** experiments. We used the **ACM** dataset for the **ACM** experiments and the **DBLP** dataset for the **DBLP** experiments. We used the **ACM** dataset for the **ACM** experiments and the **DBLP** dataset for the **DBLP** experiments.

#### Pairwise Precision

$$= \frac{\# \text{Pairs Correctly Predicted To Same Author}}{\# \text{Total Pairs Predicted To Same Author}}$$

#### Pairwise Recall

$$= \frac{\# \text{Pairs Correctly Predicted To Same Author}}{\# \text{Total Pairs To Same Author}}$$

$$\text{Pairwise } F_1 = \frac{2 \times \text{Pairwise Precision} \times \text{Pairwise Recall}}{\text{Pairwise Precision} + \text{Pairwise Recall}}$$

Our algorithm estimates  $K$  (see Algorithm 3) by iteratively finding the best two-sub-clusters model  $M_2$  for each cluster  $C$  in  $C^{(i)}$ . We calculate  $BIC(M_2)$  and  $BIC(M_1)$  for each cluster  $C$ . If  $BIC(M_2) > BIC(M_1)$ , we split cluster  $C$  into two sub-clusters  $C^{(i+1)} = \{C_1, C_2\}$ . We calculate  $BIC$  for the obtained new model, and repeat the process until no more splits are found.

1. //a e e. /d a b a



TABLE 4  
Results of Name Disambiguation (Percent)

Person Name	K-means			HAC			SOM			SACluster			CONSTRAINT			Our Approach (Fixed $K$ )		
	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1
Cheng Chang	89.47	68.00	77.27	100.0	100.0	100.0	76.30	65.42	70.44	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Wen Gao	96.23	49.78	65.62	96.60	62.64	76.38	96.52	74.14	83.68	73.52	98.27	84.11	99.29	98.59	98.94	99.29	98.59	98.94
Yi Li	13.91	39.32	20.51	86.64	95.12	90.68	43.57	32.72	37.41	77.42	84.21	80.67	70.91	97.50	87.11	70.91	97.50	87.11
Jie Tang	95.38	72.09	82.12	100.0	100.0	100.0	84.92	70.65	77.13	90.14	82.04	85.90	100.0	100.0	100.0	100.0	100.0	100.0
Gang Wu	28.41	20.49	23.81	97.54	97.54	97.54	24.79	31.28	27.66	43.66	87.32	58.22	71.86	98.36	83.05	81.62	98.36	83.05
Jing Zhang	7.88	26.03	12.10	85.00	69.86	76.69	38.76	64.23	48.35	72.00	86.75	78.69	83.91	100.0	91.25	83.91	100.0	91.25
Kuan Zhang	60.00	60.00	60.00	100.0	100.0	100.0	82.50	70.30	75.85	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Hui Fang	60.87	90.32	72.73	100.0	100.0	100.0	40.60	80.60	54.00	92.21	54.20	68.27	100.0	100.0	100.0	100.0	100.0	100.0
Wu Yan	72.25	55.00	62.50	72.25	55.00	62.50	72.25	55.00	62.50	72.25	55.00	62.50	72.25	55.00	62.50	72.25	55.00	62.50
Lei Wang	11.98	72.87	15.48	68.45	41.72	53.38	72.52	57.34	63.22	44.40	75.59	55.04	93.58	92.59	92.08	93.58	92.59	92.08
Rakesh Kumar	68.82	91.28	78.47	63.36	92.41	75.18	62.83	90.17	74.06	80.98	86.28	86.70	83.25	99.31	83.25	99.31	83.25	99.31
Michael Wagner	7.66	52.32	14.88	18.35	60.26	28.15	52.18	48.39	49.11	42.20	64.04	50.87	26.25	77.18	39.25	26.25	77.18	39.25
Bing Liu	33.10	31.73	32.42	64.88	43.10	53.22	76.80	72.60	74.64	90.21	63.63	40.63	85.72	98.83	90.55	85.72	98.83	90.55
Jim Smith	62.59	44.16	51.78	92.43	86.80	89.53	43.10	40.50	41.76	83.14	80.87	81.99	70.91	97.50	82.11	70.91	97.50	82.11
Wei Wang	11.03	10.29	10.61	8.70	10.00	9.30	10.50	10.50	10.50	12.00	66.23	20.35	33.60	94.26	44.88	33.60	94.26	44.88

fea e f eaç d; f c fee ce, e def e a e e e b e f e b e a d a d a bab ;  
 fea e a d e a e e c fee ce a e; f a , CONSTRAINT [51], a c a -ba ed c e a  
 e ea e e a a a e e, a , a f a e d a b a . F fa c a , 1) a  
 , ea , a , a d def e a fea e f eaç ba e e e, d a d, ec a ed e, d , e be K  
 a , a d, e a e b a ( dca e e ce ); f eaç a , a e ea , eac a e be ; ,  
 , ef ca , ea def e e fea e a d, e e ef a ce e e b d f e e, d; a d  
 a e e a , e de f, ec ed a e. I add , e 2) e d e e feedbac (ea ,  $r_4$ )  
 c de ed , e ba e e e, d. T ef e ba ed e e e (a , eba e e ca e, e e feedbac ).  
 , eaç ca a ea ec e (HAC) a f 5.2 Experimental Results  
 c a a d e a eaç e e , e ed a b a 5.2.1 Results  
 a [39], e a e fea e def a def ed We c d ced d a b a e e e f a e  
 ab e. T e e ba ed SAC e [52], ç e e a ed eaç f, ea , a e , e da a e. Tab e 4  
 a , e de a a , K c e b b , e e . I ca be ee , a , e, d cea  
 c a a da b e f a a ca ed eaç de. ef , eba e e e, d f a e d a b a  
 F fa c a , SAC e, e , ed , e a e (+32.77% e K-Mea , +13.28% e HAC, +33.21% e  
 a b e fea e def ed a aç a d, e a e SOM, +17.57 e SAC e, a d +10.18% e CON-  
 ea , f a . T e d ffe e ce , a SAC - STRAIN T b a e a e F<sub>1</sub> c e).  
 e d e d ffe e a e, e e f d ffe e ea , ; T e ba e e e, d ffe f d ad a a e :  
 , e c de a ea , a , e a e 1) , e ca a e ad a a e f ea , be ee  
 SAC e [52]. a e a d 2) , e e a f ed d a ce ea e.  
 We f , e c a ed e, d e A , , SAC e c de , e ea , be ee  
 e, d f a e d a b a : DISTINCT [49], a de , c a e , e ea , f a a  
 c b a e, d ba ed a ea e : e

fea e f eaç d; f c fee ce, e def e a e e e b e f e b e a d a d a bab ;  
 fea e a d e a e e c fee ce a e; f a , CONSTRAINT [51], a c a -ba ed c e a  
 e ea e e a a a e e, a , a f a e d a b a . F fa c a , 1) a  
 , ea , a , a d def e a fea e f eaç ba e e e, d a d, ec a ed e, d , e be K  
 a , a d, e a e b a ( dca e e ce ); f eaç a , a e ea , eac a e be ; ,  
 , ef ca , ea def e e fea e a d, e e ef a ce e e b d f e e, d; a d  
 a e e a , e de f, ec ed a e. I add , e 2) e d e e feedbac (ea ,  $r_4$ )  
 c de ed , e ba e e e, d. T ef e ba ed e e e (a , eba e e ca e, e e feedbac ).

## 5.2 Experimental Results

### 5.2.1 Results

We c d ced d a b a e e e f a e e a ed eaç f, ea , a e , e da a e. Tab e 4  
 , e e . I ca be ee , a , e, d cea  
 ef , eba e e e, d f a e d a b a  
 (+32.77% e K-Mea , +13.28% e HAC, +33.21% e  
 SOM, +17.57 e SAC e, a d +10.18% e CON-  
 STRAIN T b a e a e F<sub>1</sub> c e).

T e ba e e e, d ffe f d ad a a e :  
 1) , e ca a e ad a a e f ea , be ee  
 a e a d 2) , e e a f ed d a ce ea e.  
 A , , SAC e c de , e ea , be ee  
 de , c a e , e ea , f a a

TABLE 5  
Results of Our Approach with Different Settings

Method	Precision	Recall	F1-Measure
Our Approach (Auto K)	83.01	79.54	80.05
Our Approach (w/o auto K)	90.13	88.26	88.80
Our Approach (w/o relation)	67.05	50.59	55.95

f ed d a c e f c , , , ca e c d e c b e , e  
c e a b e e e , e a e a e . O f a e  
d e c d e , e c e a a , e d e d e c e  
b e e e a e e , a d e a e e e d  
a , e a e a f c b e e e a e .  
We c d c e d e , e e . T e p a e a e  
ç a e , a 0.01, d c a , a e e e  
b a a ç a e a c a f c a .  
Table 6 , e e f a a c e a f e  
be  $K$  ( e b e , e d b a c e , e a c a  
be ). We e e , a e e a e d b e b  
a a ç a e c e , e a c a b e . Table 5 f , e  
e a e a e e f a a ç d f f e e  
e , e e N / a K e e e e e f  
a a ç , a e d e f e d c e b e K a d N /  
e a , e e e e e f a a ç  
e a , ( e . , e e a e d e f e a e f c  
 $f_k(y_i, y_j)$  b e e ). We e e , a e e a  
a a a a ç . W , e e a , e  
e f a c e f a a ç d , a (-23.08 e c e  
b  $F_1$  c e ). T c f , a a d e , ç c a  
c a e d e d e c e b e e e a e d e  
d e f a c e .  
We a e d X- e a f d e b e f e e K .  
We a e d e b e a l a d a  
b e a  $n$  , e a e e a a . We  
f d , a X- e a f a f d e a c a b e . I  
a a e c e e c e N L , 2. T e  
e a b e , a X- e a c a a e e f e  
e a b e e e a e .

TABLE 6  
Result of Automatically Discovered Person Number

Person Name	Actual Number	Auto Number	Person Name	Actual Number	Auto Number
Cheng Chang	3	3	Dimitry Pavlov	2	1
Wen Gao	4	5	David Jensen	3	6
Yi Li	21	13	David Brown	7	9
Jie Tang	2	2	David C. Wilson	5	5
Gang Wu	16	12	George Miller	2	6
Jing Zhang	25	16	James H. Anderson	2	7
Kuo Zhang	2	2	James Johnson	3	3
Hui Fang	3	3	John Miller	2	5
Bin Yu	12	10	Joseph Miller	2	3
Lei Wang	40	22	Paul Jones	3	5
Rakesh Kumar	5	5	Richard Taylor	10	14
Michael Wagner	10	11	Robert Fisher	4	7
Bing Liu	11	12	Robert Moore	3	6
Jim Smith	5	5	Robert Williams	2	5
Wei Wang	90	22	William Cohen	2	9
Ajay Gupta	4	6	Charles Smith	4	

TABLE 7  
Comparison with DISTINCT

Person Name	DISTINCT			Our Approach		
	Prec.	Rec.	F1	Prec.	Rec.	F1
Cheng Chang	55.07	44.19	49.03	100.00	100.00	100.00
Wen Gao	92.07	98.68	95.26	90.20	98.50	94.04
Jie Tang	79.36	93.37	85.80	100.00	100.00	100.00
Jing Zhang	100.00	75.56	86.08	83.91	100.00	91.25
Kuo Zhang	78.57	84.78	81.56	100.00	100.00	100.00
David Jensen	85.69	100.00	92.29	83.83	68.46	75.37
David Brown	69.77	74.99	72.29	89.32	91.45	90.37
David C. Wilson	87.10	90.00	88.53	94.33	67.30	78.55
Richard Taylor	68.35	63.11	65.63	94.33	79.00	86.41
Charles Smith	78.42	76.67	77.54	100.00	100.00	100.00
Hui Fang	88.60	95.00	91.69	100.00	100.00	100.00
Rakesh Kumar	92.90	96.80	94.81	99.14	96.00	98.01
Michael Wagner	72.30	75.40	73.82	85.69	82.31	83.97
Bing Liu	78.30	95.70	86.13	88.25	86.00	87.36
Jim Smith	86.30	90.40	88.30	96.37	93.80	95.07
Lei Wang	80.80	89.60	84.97	89.17	88.00	89.05
Bin Yu	68.90	77.80	73.08	95.27	72.63	82.42
Wei Wang	78.60	78.30	78.45	85.19	83.65	84.14
Ajay Gupta	98.70	92.30	95.39	97.78	99.65	97.11
Avg.	81.04	83.82	82.14	93.78	89.80	91.48

We c a e d a a ç , DISTINCT [49]. We  
e d e a e , a e e e d b , [49] a d  
e e e f c a . We c d c e d e e e -  
e e e d a a e , ç a e e e e f d a e d  
[49]. F e a e , e a e 109 a e f 'N e W a ,  
a d 33 a e f 'N S , e [49] , e b e  
a e 55 a d 19. I a d d , e d c d e , e  
P c e e d E d e a . Table 7 , e c a  
e . We e e , a a e a e , d c e a e -  
f DISTINCT (+8.34% b  $F_1$ ). M e e , a a ç  
a , e a d a e a c a a c a f d e b e  
K , e e a DISTINCT e b e e d b e e d  
b , e e . T e e a e d DISTINCT a d  
a a ç a e d f f e e . DISTINCT a c d e , e  
a , - a e a d a e - c f e e c e e a , a d d e  
d e c c d e , e C A , a d C P b V e e e a ,  
a , e e e a c a b e d e d f , e a e -  
c f e e c e a d a , - a e e a .

## 5.2.2 Efficiency Performance

We e a a e d , e e f f c e c e f a c e f a a ç  
f , e 32 a , a e a d e c e . I e  
C e D c e (1.6 GH ). Table 8 , e CPU  
e e d f a , e a e d f f e e a , . We  
a , b , e a 100 a e a d  
e a e a e e f 100 a d a e . F a ,  
a e , a e a , e e , a 1 e c d . T e a  
e f a a , a e a ç , e .

TABLE 8  
Comparison of Efficiency Performance (Seconds)

Person Name	K-means	K-Means	HAC	SA Cluster	DISTINCT	Our Approach
Wen Gao	4.8	3.1	12.9	30.4	36.0	20.3
Lei Wang	3.7	2.4	6.8	4.1	12.1	4.6
Bing Liu	1.6	1.9	4.2	5.4	1.1	5.8
Wei Wang	28.7	5.1	73.1	46.9	83.3	100.2
Robert Fisher	2.8	1.3	5.6	0.2	0.2	0.8
William Cohen	0.8	1.2	3.0	0.06	0.6	0.9
Average over 100	0.52	0.26	1.14	0.96	0.87	1.42

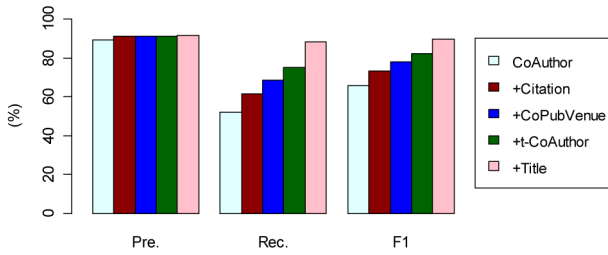


Fig. 3. Contribution of relationships.

### 5.2.3 Feature Contribution Analysis

We e a ed , e c b f , e def ed fea e ( c d ed ea d defea e ) f a ed a b a . Sec f ca , e f a , e d d a fea e b , e e f a ce , e add , e fea e eb e , e de f , e d a b a e . I a c a , e f e C A , f ed b add C a , a d , e C P bVe , Pa e T e . I eaç e , ee a ae , e ef a ce f e , d . F . 3 , e a e a e Pec , a e a e Reca , a d a e a F1- c e f e , d d ffe e fea ec b a . A eaç e , e b e ed e . We ca a ee , a f , e fea e ( e ce C A ) a c b e , e e e f e ca , e , e e e e ec ed .

### 5.2.4 Distribution Analysis

We a e f ad b a a ad e ed c e , d [10]. We f d , a , e fea ed b - f a a e ca be ca ca e ed , e f ce a : 1) b ca f d ffe e e a e cea e a a ed (NH Fa , ). Na e d a b a , d f d a ca be ed e e b a a aç a d , e be K ca a bef d acc a e ; 2) b ca a e ed e , e b a d a a , e f , e a e (e . , NB L . ) ; a aç ca aç e e a F1 c e f 87.36 e ce a d , ed c e ed be K c e , e ac a be ; a d 3) b ca f d ffe e a , a e ed (e . , N Z a , ). O e , d ca b a a ef a ce f 91.25 e ce . H e e , d be d ff c acc a e f d , e be K . F e a e , e be f d b a aç f N Z a , 14 , b , e c ec be , d be 25 . F a de a ed a a , ea e efe [41].

### 5.2.5 Application Experiments

We a ed , e a ed a b a , e e e f d , ç de f e , e e e e e e e e ce . I a c a , ee a a ed e e f d a d , a ed a b a . Sec f ca , e e e ed 12 f e e e e f , e e f A e M e , a d , ed e e a ce d e [6] e , e a d e ce a e a da a e f e a a . I e e ed eade a e efe ed [51], [40] f de a f , e e e e a e . We c d c ed e a a e f P@5, P@10, P@20, P@30, R- ec, ea a e a e ec (MAP), bpref, a d ea ec ca a (MRR). F . 4 , e e f e e f d . I F . 4, EF e e e e e f d a ed a b a b e , da d EF- NA e e e e e f d , a ed a b a . We ee , a cea e e ca be ba ed b , e ed a ed a b a a aç .

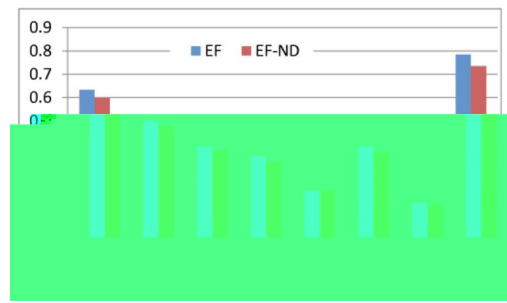


Fig. 4. Performances of expert finding.

## 5.3 Online System

T f , e de a e , e effec e e f , e ed a aç , e , a e a ed , ed a b a , e , d , e A e e e . F . 5 , a a , f , e d a b a e . T e e eaç e f N e Ta , a d , e e e , ee d ffe e e , e f , e a e a d be , e de a ed f e f a f eaç e . T e e , d a ff e de a d fa e e a ead e e a e , ed a b a e f e , a 10,000 e a e . Pea e e , a a ec . V , de ec , e e ç a e .

## 6 DISCUSSION

### 6.1 Connections with Previous Work

We a a e , e c ec f f a e , e e a e e d a b a / c . e .  
**Connection with K-means:** O f a e ca de c be ea , be ee da a , e ea K- ea [27] ca . I e e ce , f a e e ed e e a f c de , e ea . B e , e ed e e a f c f (8), e , a e

$$L_{\max} = \sum_{x_i \in X, l} \alpha_l K(x_i, \mu_l) - \log Z. \quad (19)$$

B f , e e , e e ,  $\alpha_l$  f eaç a f c , e b a a a e K- ea c e a .  
**Connection with X-means:** X- ea [33] ed d a ca f d , e c e be K . I a e BIC f de e ec . H e e , a de d ffe a e f X- ea , e e ec ce a d , e c e a , a e a d ffe e . T e de e ec

Fig. 5. Name disambiguation system (<http://arnetminer.org>).

e, d f a e a a X- ea f e  
e c de e bab  $P(Y)$  f, .e.,  
de e de ce be ee daa . E ce f  
de eec , X- ea ef a K- ea .

**Connection with the constraint-based disambiguation method:** I c a -ba ed c e ,e., [2], e e ca

c a de e c e ce . I  
e, a ea ed a e d a b a a d  
b a ed e e [51], [41]. T e a c a  
c de - a dca - . M - ea , a  
da a be ed ec e a dca -  
ea da a be ed d ffe e  
c e . We ca ada f a e a c a -ba ed  
c e b edef , eed e e a f c .

**Connection with disambiguation using spectral graph clustering:** S ec a a , c e [12] a a f d

b a , c f ea , be ee da a  
K- a ec a a , c e a , a bee  
e ed f a e d a b a [18]. We ca e a  
e a , e ed da a a f , e e e a ed  
d ffe e c e (.e.,  $I(i \neq j)$ ) , e bec ef c .  
T e , f a e ca ada , c e b e  
e ec d a f (8)

$$L_{\min} = - \sum_{(x_i, x_j) \in E, R, k} K(x_i, x_j) r_k(x_i, x_j) + \log Z. \quad (20)$$

I e e ce, e bec ef c ea , a e  
e, e e a e bab e , e HMRf a d  
f c , e de e de ce be ee a e .  
C a , e e f a e  
ffe e e a ad a a e : 1) I ad a e , d , a -  
e f a e a e de e de , , ca a e  
ad a a e f ea , be ee a e . 2) T e -  
ed f a e ca be ea e e ded e - e -  
ed ea b e feedbac . 3) O  
f a e ca be e ed a a e e a f a e f  
e e a e ed e , d .

## 7 CONCLUSION AND FUTURE WORK

I , a e, e a e e a ed , e be f a e  
d a b a . We a ef a ed , e be a  
fed f a e ad ed a e e a ed bab -  
c de , e be . We a ed ef ed a d a b a -  
bec ef c f , e be a d a e ed a  
- e a a e e e a a . We a e a  
e ed a d a ca a c f e a , e be f  
e e K. E e e a e d ca e , a e ed  
e, d fca ef , e ba e e e , d .  
W e a ed e e f d , cea e e (+2%)  
ca be b a ed.

A , e e e , d be e e e a e  
a e e f e e f a f a e  
d a b a , a e a b be e e  
e e . M e e , a e e d , c  
de e LDA ca e a ed a b a .

## ACKNOWLEDGMENTS

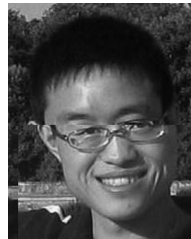
T e a , d e , a H C e f d  
e c e c de f SAC e a d X a Y f d  
e c e c de f DISTINCT f , ec a e e -  
e . T e a , a P f . P . Y f , a abe

. Je Ta ed b , e  
Na a Sc e ce F da f C a (N . 61073073), e  
C e e Na a Ke F da Re ea c (N . 60933013,  
N . 61035004), a d a S ec a F d f FSSP.

## REFERENCES

- [1] H. A a e, 'NA Ne L a , e Sa ca M de Ide f ca-  
IEEE Trans. Automatic Control, . AC-19, . 6,  
. 716-723, Dec. 1974.
- [2] S. Ba , M. B e , a d R.J. M e , 'NA P bab c Fa e-  
f Se -S e ed C e , Proc. ACM SIGKDD Int'l  
Conf. Knowledge Discovery and Data Mining (SIGKDD '04), . 59-  
68, 2004.
- [3] R. Be e a a d A. McCa , 'ND a b a Web A ea -  
a ce f Pe e a S ca Ne , Proc. Int'l Conf. World Wide  
Web (WWW '05), . 463-470, 2005.
- [4] O. Be e , H. Ga ca-M a , D. Me e a , Q. S , S.E.  
W a , a d J. W d , 'NS : A Ge e c A a c E  
Re , The VLDB J., . 18, . 255-276, 2008.
- [5] I. B a a , a a d L. Ge , 'NC ec e E Re  
Re a a Da a , ACM Trans. Knowledge Discovery from Data,  
. 1, a c e 5, 2007.
- [6] C. B c e a d E.M. V , ee , 'NRe e a E a a  
I c ee I f a , Proc. Ann. Int'l ACM SIGIR Conf.  
Research and Development in Information Retrieval (SIGIR '04),  
. 25-32, 2004.
- [7] Z. C e , D.V. Ka a , a d S. Me a , 'NAda e G a , ca  
A a c E Re , Proc. Seventh ACM/IEEE-CS Joint  
Conf. Digital Libraries (JCDL '07), . 204-213, 2007.
- [8] Z. C e , D.V. Ka a , a d S. Me a , 'NE C e  
A a f C b M e E Re S e ,  
Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '09),  
. 207-218, 2009.
- [9] D. C , R. Ca a a , a d A. McCa , 'NSe - e ed  
C e U e Feedbac , Tec ca Re TR2003-1892,  
C e U , 2003.
- [10] D. Ca , X. He, a d J. Ha , 'NS ec a Re e f D e a  
Red c , ec ca e , 2856, UIUC 2004.
- [11] P.T. Da , D.K. E , a d J.L. Ka a , 'NMe , d f Pec e  
Na ed E Ma c D a C ec , Proc. ACM/IEEE-  
CS Joint Conf. Digital Libraries (JCDL '03), . 125, 2003.
- [12] C. D , 'NA T a S ec a C e , Proc. Int'l Conf.  
Machine Learning (ICML '04), 2004.
- [13] M. E e , R. Ge, B.J. Ga , Z. H , a d B. Be -M , e, 'N  
C e A a f A b e Da a d Re a Da a : T e  
C ec ed K-Ce e P be , Proc. SIAM Conf. Data Mining  
(SDM '06), 2006.
- [14] S. Ge a a d D. Ge a , 'NS c a c Re a a , Gbb D b -  
a d , e Ba e a Re a f I a e , IEEE Trans. Pattern  
Analysis and Machine Intelligence, . PAMI-6, . 6, . 721-742,  
N . 1984.
- [15] Z. G a a a a d M.I. J da , 'NFac a Hdde Ma  
M de , Machine Learning, . 29, . 245-273, 1997.
- [16] J. Ha e e a d P. C ff d, 'NMa Fed F e G a ,  
a d La ce , U b , ed a c , 1971.
- [17] H. Ha , L. Ge , H. Z a , C. L , a d K. T , 'N  
S e ed Lea A a c e f Na e D a b a ,  
A , C a , Proc. ACM/IEEE Joint Conf. Digital Libraries  
(JCDL '04), . 296-305, 2004.
- [18] H. Ha , H. Z a , a d C.L. Ge , 'NNa e D a b a A ,  
C a U a K-Wa S ec a C e Me , d , Proc. ACM/  
IEEE Joint Conf. Digital Libraries (JCDL '05), . 334-343, 2005.
- [19] G.E. H , 'NT a P d c f E e b M  
C a e D e e ce , J. Neural Computation, . 14, . 1771-  
1800, 2002.
- [20] L. J a , J. Wa , N. A , S. Wa , J. Z a , a d L. L , 'NGRAPE: A  
G a , -Ba ed Fa e f D a b a Pe e A ea a ce  
Web Sea c , Proc. Int'l Conf. Data Mining (ICDM '09), . 199-  
208, 2009.
- [21] M.I. J da , Z. G a a a , T. Jaa a , a d L. Sa , 'NA  
I d c Va a a Me , d f G a , ca M de ,  
Learning in Graphical Models, . 37, . 105-161, 1999.
- [22] R. Ka a d L. Wa e a , 'NA Refe e ce Ba e a T e f Ne ed  
H , e e a d I Re a , e Sc a C e , J. Am.  
Statistical Assoc., . 90, . 773-795, 1995.

- [23] L. Ka f a a d P. R e e , *Finding Groups in Data: An Introduction to Cluster Analysis*. We e , 1990.
- [24] R. K d e a a d J.L. S e , *Markov Random Fields and Their Applications*. A . Ma . S c ., 1980.
- [25] H. K c , S. G e a , a d A. K e a a , 'N d d e Ma Ra d F e d , *J. Annals of Applied Probability*, . 5, . 3, . 577-602, 1995.
- [26] X. L , P. M e , D. R , 'N d e f c a a d T a c f A b Na e : D c a e a d G e e a e A a c e , *Proc. 19th Nat'l Conf. Artificial Intelligence (AAAI '04)*, . 419-424, 2004.
- [27] J. MacQ e e , 'N S e M e , d f C a f c a a d A a f M a a e O b e a , *Proc. Fifth Berkeley Symp. Math. Statistics and Probability*, 1967.
- [28] D.M. McRae-S e e a d N.R. S a d b , 'N A b , e S a e A , : A K T e A , a C a U G a , A a c Na e D a b a , *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '06)*, . 53-54, 2006.
- [29] E. M , W.W. C e , a d A.Y. N , 'N C e a S e a c a d Na e D a b a E a U G a , *Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '06)*, . 27-34, 2006.
- [30] K.P. M , Y. W e , a d M.I. J d a , 'N L B e e f P a a f A a e I f e e c e A E c a S d , *Proc. Conf. Uncertainty in Artificial Intelligence (UAI '99)*, . 467-475, 1999.
- [31] M.E.J. N e a a d M. G a , 'N F d a d E a a C S c e N e , *Physical Rev. E*, . 69, . 026113, 2004.
- [32] B. O a d D. L e e , 'N S c a b e N a e D a b a U M - L e e G a , P a , *Proc. SIAM Int'l Conf. Data Mining (SDM '07)*, 2007.
- [33] D. P e e a d A. M e , 'N K-Me a : E e d K-Me a E f f e E a f , e N b e f C e , *Proc. Int'l Conf. Machine Learning (ICML '00)*, 2000.
- [34] J. R a e , 'N A U e a P f I e e a d E a b M D e c L e , *J. Annals of Statistics*, . 11, . 2, . 416-431, 1983.
- [35] J. S a d J. M a , 'N N a e d C a d I a e S e e a , *IEEE Trans. Trans. Pattern Analysis and Machine Intelligence*, . 22, . 8, . 888-905, A . 2000.
- [36] L. S , B. L , a d W. M e , 'N A L a e T c M d e f C e e E R e , *Proc. IEEE Int'l Conf. Data Eng. (ICDE '09)*, . 880-891, 2009.
- [37] Y. S , J. H a , I.G. C c , J. L , a d C.L. G e , 'N e f f c e T c-Ba e d U e e d N a e D a b a , *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '07)*, . 342-351, 2007.
- [38] Y. S , Y. Y , a d J. H a , 'N R a -B a e d C e f H e e e I f a N e S a N e S c e - a , *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '09)*, 2009.
- [39] Y.F. T a , M. K a , a d D. L e e , 'N S e a c E e D e A , D a b a , *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '06)*, . 314-315, 2006.
- [40] J. T a , J. Z a , L. Y a , J. L , L. Z a , a d Z. S , 'N A e M e : E a c a d M f A c a d e c S c a N e , *Proc. 14th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '08)*, 2008.
- [41] J. T a , L. Y a , D. Z a , a d J. Z a , 'N A C b a A a c Web U e P f , *ACM Trans. Knowledge Discovery from Data*, . 5, a c e 2, Dec. 2010.
- [42] Y. T a , R.A. H a , a d J.M. P a e , 'N e f f c e A e a f G a , S a a , *Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '08)*, . 567-580, 2008.
- [43] J. V e a a d E. A , e , 'N C e f , e S e f-O a M a , *IEEE Trans. Neural Network*, . 11, . 3, . 586-600, Ma 2000.
- [44] M. W e a d G.E. H , 'N A N e L e a A f M e a F e d B a M a c e , *Proc. Int'l Conf. Artificial Neural Networks (ICANN '01)*, . 351-357, 2001.
- [45] M. W e a d K. K , a a , 'N B a e a K-Me a a a N M a a -E e c a , A , *Proc. SIAM Int'l Conf. Data Mining (SDM '06)*, . 472-476, 2006.
- [46] S.E. W a , D. M e e a , G. K , a , M. T e b a d , a d H. G a c a-M a , 'N E R e I e a e B c , *Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '09)*, . 219-232, 2009.
- [47] S.E. W a , O. B e e , a d H. G a c a-M a , 'N G e e c E R e N e a e R e , *The VLDB J.*, . 18, . 6, . 1261-1277, 2009.
- [48] X. X , N. Y , Z. F e , a d T.A.J. S c e e , 'N S c a C e A f N e , *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '07)*, . 824-833, 2007.
- [49] X. Y , J. H a , a d P.S. Y , 'N O b e c D c : D O b e c I d e c a N a e , *Proc. Int'l Conf. Data Eng. (ICDE '07)*, . 1242-1246, 2007.
- [50] H. Y , W. K , V. H a a , a d J. W b , 'N A L a e S c a e , C -B a e d A a c f A a c a D a b a B -e d c a A b b e a , *ACM Trans. Information Systems*, . 24, . 3, . 380-404, 2006.
- [51] D. Z a , J. T a , J. L , a d K. W a , 'N A C a -B a e d P b a b c F a e f N a e D a b a , *Proc. ACM Conf. Information and Knowledge Management (CIKM '07)*, . 1019-1022, 2007.
- [52] Y. Z , H. C e , a d J.X. Y , 'N G a , C e B a e d S c a A b e S a e , *Proc. VLDB Endowment*, . 2, . 1, . 718-729, 2009.



**Jie Tang** is an associate professor at Tsinghua University. His research interests are social network analysis, data mining, and semantic web.



**A.C.M. Fong** is a professor in the School of Computing and Mathematical Sciences, Auckland University of Technology. He has published widely in the areas of data mining and communications.



**Bo Wang** is currently working toward the PhD degree from Nanjing University of Aeronautics and Astronautics. His research interests include transfer learning and information network analysis.



**Jing Zhang** received the MS degree from Tsinghua University in 2008. Her research interests include information retrieval and text mining.

► For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).