

# Study on Paradigm and Its Normalization with Coexistence of XML Functional Dependence and Multivalued Dependence

Lijun Cao

Hebei Normal University of Science & Technology

Xiyin Liu, Jing Cao, Aiyong Liu, Zhongping Zhang

Hebei Normal University of Science & Technology, China

College of Information Science and Engineering Yanshan University, China

**Abstract**—This paper proposes the paradigm and the corresponding normalization algorithms with coexistence of eXtensible Markup Language (XML) functional dependence and XML values dependence. The concepts of redundancy and keys are firstly described in this paper, while the paradigm with coexistence of XML functional dependence and XML values dependence is defined later. Then, the non-redundancy judgment theorem is also presented. Furthermore, the XML document tree normalization algorithm is introduced, while the capability to be terminated, validity, and computation complexity are analyzed and proved. Finally, the effectiveness of the algorithm is demonstrated by the experiments.

**Index Terms**—normalization; redundancy; paradigm; functional dependence; multi-value dependence; normalization rules

## I. INTRODUCTION

While XML is widely used in more and more fields, a well-performed eXtensible Markup Language (XML) document with more Document Type Definitions (DTDs) is necessary to reduce the data redundancy. Model normalization theory is the theoretical basis to design a well-performed database model. Similar to relation databases, if the DTD design pattern is not suitable enough, abnormal data operation will take place. To fix this problem, normalization study is a good method to eliminate data redundancy and abnormal operation. Therefore, the paradigm theory is one of the hot topics in the research of database [1].

For the cases with coexistence of XML functional dependence and multivalued dependence, this paper proposes a paradigm and its normalization algorithm based on their logical implication and membership. After studying the relation database normalization theorem and characteristics of XML database, we also define 3XNF (XML Normal Form) and 4XNF, and further give normalization rules and a theorem on non-redundancy judgment. The normalization algorithm

is finally designed while its effectiveness is then proved by the experiments.

## II PARADIGM IN XML FUNCTIONAL DEPENDENCE

In relation databases, the usage of paradigm can reduce data redundancy and decrease the number of abnormal operations. The concept of paradigm is also utilized in XML. People can discuss whether a DTD is well-designed using a paradigm, and then apply normalization operation on the DTDs which do not meet the paradigm.

### A. Effective Change and Redundancy

While XML is widely used nowadays, the design of a good model is becoming an important research topic [2]. Data redundancy not only wastes a lot of space for storage, but also increases the cost of data transfer and data operation. Therefore, a model with a good design pattern is needed to reduce data redundancy.

In relation databases, data redundancy means the repeated saving of the same data for many times. In a relation mode which satisfies dependence set  $\Sigma$ , if an attribute value can be obtained according to the dependence set as well as other attribute values, this attribute value is believed to be redundant. Particularly, there is a relationship instance based on the model, which satisfies the dependence set  $\Sigma$ . If an instance will not satisfy some dependencies in set  $\Sigma$  after a value in this instance has been changed, the relation model is called to be redundant.

When borrowing the concept of redundancy theory in relation database, we present the definitions of effective change and redundancy as follows.

**Definition 1 (Effective change)** For a given DTD:  $D$ ,  $T$  is a complete tree satisfying  $D$ . Let  $v$  be an arbitrary node in  $T$ .  $v$  is changed into  $v'$ , where the identifier  $v'$  is different from other identifiers in  $T$ , and  $\text{val}(v) \neq \text{val}(v')$ , a new tree  $T'$  is generated after this changing. If  $T'$  still satisfies  $D$ , then  $T$  is called as an effective change.

**Definition 2 (Redundancy)** For a given DTD:  $D$ ,  $\Sigma$  is a collection set of XFD and XMVD (XML Multi-valued

Dependency), where the paths in XFD and XMVD are all collected in the set Paths (D). In a complete tree T satisfying D and  $\Sigma$ , a node v is changed into v'. If this changing will generate a new tree T', which violates  $\Sigma$ , D is said to be redundant.

*B. The Third Paradigm in XML (3XNF)*

In relation databases, BCNF is a paradigm with the highest normalization, which only considers the cases with functional dependence [3], while 3NF is a paradigm with the highest normalization to keep the functional dependence. Borrowing the related theory in relation data, we present the formal definition of the third paradigm in XML (3XNF) as below.

Definition 3 (The third paradigm in XML) For a DTD: D,  $\Sigma$  is the smallest dependance set in XFD. If for each abnormal XFD with  $p_1, \dots, p_k \rightarrow q_i$ , there exists XFD with  $p_1, \dots, p_k \rightarrow pparent(q)$ , where,  $last(q_i) \in E2 \cup A$ , then D is called to be the third paradigm in XML.

The definition of XML abnormal functional dependence which generates redundancy is given as follows.

Definition 4 (XML abnormal functional dependence) Assume there is a given DTD: D and a smallest dependance set containing XFD. For each abnormal XFD with  $p_1, \dots, p_k \rightarrow q_i$ , if this XFD does not satisfy the definition of XNF, i.e.  $p_1, \dots, p_k \rightarrow q_i$  is established in D, but  $p_1, \dots, p_k \rightarrow pparent(q)$  is not, where  $last(q_i) \in E2 \cup A$ , then XFD:  $p_1, \dots, p_k \rightarrow q_i$  is called to be abnormal functional dependence.

Theorem 1: Assume in a given DTD D, T is a complete tree satisfying D,  $\Sigma$  is the XFD collection set, and T satisfies  $\Sigma$ . If D satisfies the 3XNF model, then T is called to be with non-redundancy.

Proof: This theorem is proved by reduction to absurdity. That is: If T is redundant, then D does not satisfy 3XNF mode.

Assume there is redundancy in T, following the definition of redundancy, we have: there exists a node v in T with  $v \in N(q)$ , and abnormal XFD:  $P \rightarrow q$  is not established when v turns to v' with effective change.. According to definition 3.1, there are two different route instances  $q_1: v_1, \dots, v_n (v_n=v)$  and  $q_2: w_1, \dots, w_n (w_n=v)$  in Paths(q).

(1) When  $pparent(q_1) \neq pparent(q_2)$ , we have below conclusions we have conclusions as given below . because of  $val(v_n)=val(w_n)$  and  $pparent(q_1) \neq pparent(q_2)$ ,  $P \rightarrow parent(q)$  is not established according to the definition of XFD, and thus it does not satisfy the definition of 3XNF. While When considering that T is a complete tree satisfying D, it's known that D does not satisfy 3XNF.

(2) When  $pparent(q_1)=pparent(q_2)$ , below discussion can be established. Because  $last(q) \in E2 \cup A$ ,  $q_1$  and  $q_2$  are the same route according to definition 2.2. This conflicts with the assumption that  $q_1$  and  $q_2$  are different.

In conclusion, if D satisfies the 3XNF model, then T is called to be with non-redundancy.

*C. Design and Algorithm for XML Function Normalization*

This section will firstly give normalization rules, then discuss how to change any DTD into normalization model to meet the paradigm of 3XNF, and further present the normalization algorithm.

XML functional dependence normalization rules are given as follows.

(1) Normalization rule 1: move sub-tree with attributes and simple elements.

If the complex element vparent (last (q)) is embedded in another complex element last(P), then vparent (last (q)) is with some attributes or with simple elements, while these attributes or simple elements are closer with last (P). In order to reduce the data redundancy, we can move attributes/child elements with simple elements being last (P), to construct a new DTD D' as shown in below figure to reduce the data redundancy.

Given a DTD: D,  $\Sigma$  is a minimum function dependence set on XFD in D. Assume there exists an abnormal functional dependence  $p_1, \dots, p_k \rightarrow q$  in D, where  $last(q) \in E2 \cup A$ . If there exists a left redundancy route, i.e.  $p_1, \dots, p_i \rightarrow q$  is established while  $P' = \{ p_1, \dots, p_i \}$ , where p' is the public prefix for all routes in P', then last(q) is a child element of last(p'), i.e.  $vparent(last(q))=last(p')$ , where,  $p_j$  satisfies  $pparent(p_j)=p'$ , and  $p_j \subset q, 1 \leq j \leq i$ .

After DTD D has been normalized, paradigm can be represented as  $D' = (E', A', P', R', r)$ , where  $E' = E; A' = A; P' = P$   
 $(last(p') \rightarrow last(q)) = P(last(p') \rightarrow last(q)) \cup vparent(last(q)); P' (vparent(last(q))) = P(vparent(last(q))) - \{last(q)\};$

$R' = R$ ; other definitions of P' are same to that of P. Let XFD set of D' be  $\Sigma'$ , this kind of functional dependence is turned into  $p_1, \dots, p_k \rightarrow q$  in  $\Sigma'$ , i.e.  $\Sigma' = \Sigma - p_1, \dots, p_k \rightarrow q \cup p_1, \dots, p_k \rightarrow p' \rightarrow vparent(last(q)).last(q), k \geq 1$ .

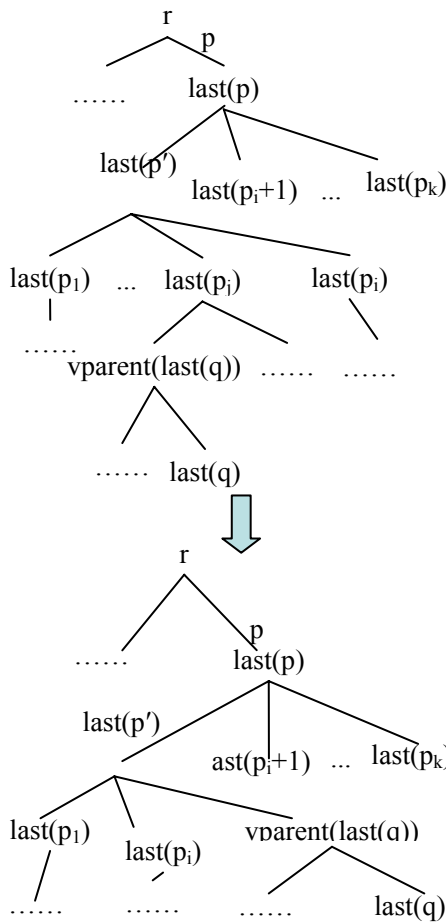


Figure 1 mobile attributes and sub-tree containing simple elements

(2) Normalization rule 2 Generation of element node

If the complex element  $vparent(last(q))$  is embedded in a complex element  $last(P)$ , attributes or simple type element in  $vparent(last(q))$  have no direct relationship with  $last(P)$ . This is the time to create new element to reduce redundancy, while the attributes of the new node are that of  $vparent(last(q))$  in the tree. Child elements are the simple attributes of  $vparent(last(q))$  in the original tree. DTD  $D'$  has been constructed as below figure to reduce the data redundancy.

Given a DTD:  $D$ ,  $\Sigma$  is a minimum function dependence set on XFD in  $D$ . Assume there exists an abnormal functional dependence  $p_1, \dots, p_k \rightarrow q$  in  $D$ , where  $last(q) \in E_2 \cup A$ . Let  $p_1 \cap \dots \cap p_k = p, ancestor(q) \cap p = s$ , if there exists  $last(p_i) \in E_2 \cup A$ , we move  $last(q)$  to make it be a child element of  $last(s)$  i.e.  $vparent(last(q)) = last(s)$ .

In the left route set  $P = \{p_1, \dots, p_k\}$ , let all routes in  $last(p_i) \in E_2 \cup A$  be added in set  $Vpath(P)$ , i.e.  $Vpath(P) = \{p_g, \dots, p_h\}$ , let all routes in  $last(p_i) \in E_1$  be added in set  $Epath(P)$ , i.e.  $Epath(P) = \{p_e, \dots, p_f\}$ .

After DTD  $D$  has been normalized, paradigm can be represented as:  $D' = (E', A', P', R', r)$ , where,  $E' = E \cup \{lab(vnew)\}$ ;  $A' = A$ ;  $P' (last(s)) = P(last(s)) \cup \{lab(vnew)\}$ ;  $P' (last(vnew)) = \{last(p_g), \dots, last(p_h), last(q)\}$ ;

$P'(vparent(last(q))) = P(vparent(last(q))) - \{last(q)\}$ ; other definitions of  $P'$  are same to that of  $P$ .  $R'(lab(vnew)) = \{vi | vi = last(pi) \in A \wedge last(pi) \in P (last(vnew))\}$ ;  $R'(vparent(last(q))) = R(vparent(last(q))) - \{last(q) | last(q) \in A\}$ ; other definition of  $P'$  and  $R'$  are same with that in  $D$ . Let the functional dependence set in  $D'$  be  $\Sigma'$ , this kind of functional dependence in  $\Sigma$  is turned into the formation of s. vnew.  $last(p_g), \dots, s. vnew. last(p_h) \rightarrow s. vnew. last(q)$ , i.e.  $\Sigma' = \Sigma - p_1, \dots, p_k \rightarrow q \cup s. vnew. last(p_g), \dots, s. vnew. last(p_h) \rightarrow s. vnew. last(q), k \geq 1$ .

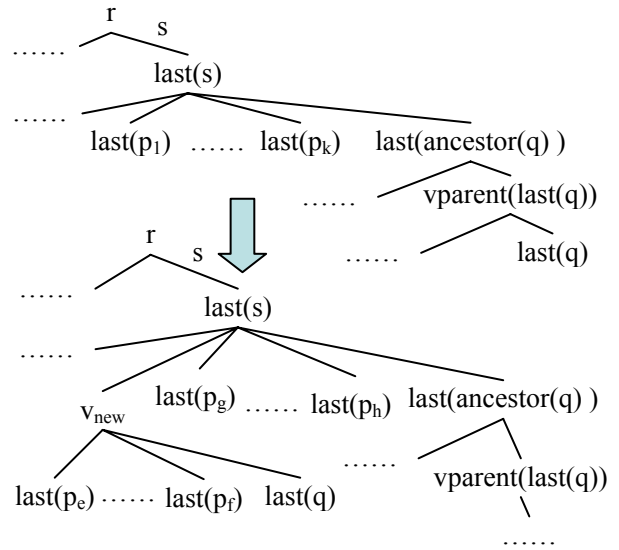


Figure 2 Generation of element node

First of all,  $D'$  is initialized as  $D$ , and  $\Sigma'$  is initialized as  $\Sigma$ . We first need to find out whether  $D'$  satisfies 3XNF. If yes, it directly returns  $D$ ; if not, it indicates that there exists abnormal functional dependence, and operation given below should be done. If the right tail node of the abnormal functional dependence is a simple element or attribute node, wherein the left path of this dependence is with public prefix of  $p$  and we find that  $\Sigma$  containing logically  $p_1, \dots, p_k \rightarrow p$  according to the membership algorithm, then normalization can be done with moving elements and following normalization rules. Otherwise, normalization will be done by generating element node normalization rules.

Algorithm 1: Algorithm of converting DTD  $D$  into DTD  $D'$ , which satisfies 3XNF (DTD-to-3XNF).

Input: DTD  $D$  and functional dependence set  $\Sigma$ ;

Output: DTD  $D'$  satisfying 3XNF

DTD-to-3XNF( $D, \Sigma$ )

Begin

(1) Initialization:  $D' := D, \Sigma' := \Sigma$ ;

(2) If  $D'$  satisfies 3XNF, turn to (6);

(3) While each abnormal XFD  $\in F$  do

(4) Case1: if  $p_1, \dots, p_k \rightarrow q \in F$  and  $p \leq q$ , where  $last(q) \in E_2 \cup A, p_1 \cap \dots \cap p_k = p$  and  $DEP\_MEMBERSHIP(\Sigma, p_1, \dots, p_k \rightarrow p)$  then use normalization 1;

(5) Case2:if  $p_1, \dots, p_k \rightarrow q$  and  $\text{last}(p_i) \in E_2 \cup A$ , where  $\text{last}(q) \in E_2 \cup A$ ,  $s = \text{ancestor}(q) \cap p_1 \cap \dots \cap p_k = p$  and  $\text{DEP\_MEMBERSHIP}(\Sigma, p_1, \dots, p_k \rightarrow p)$  then use normalization 2;

(6) Return( $D'$ ).

End

The analysis of algorithm 1 is as below.

(1) The correctness of the algorithm

The algorithm presents the corresponding solutions according to the different situations. It makes XFD satisfying its definition using not only the method of moving the attributes and simple element, but also the method of creating and copying element nodes. If the definition of XNF is satisfied, apparently the correctness of the algorithm is established.

(2) The capability to be terminated

This capability is determined by the number of the abnormal functional dependence because this number is limited, and can be terminated. Therefore, the whole algorithm can be terminated.

(3) The time complexity of the algorithm

Let  $B$  be the number of abnormal XFD in  $\Sigma$ . The algorithm's time complexity depends on the time complexity of  $\text{DEP\_MEMBERSHIP}$  algorithm and the number of execution of  $\text{DEP\_MEMBERSHIP}$ , while the  $\text{DEP\_MEMBERSHIP}$  algorithm's time complexity is  $O(n \times m^3)$ .  $\text{DEP\_MEMBERSHIP}$  execution times are determined by the number of abnormal XFD  $b$ . Therefore, the algorithm's time complexity is  $O(b \times n \times m^3)$

### III PARADIGM WITH COEXISTENCE OF XML FUNCTIONAL DEPENDENCE AND XML VALUES DEPENDENCE

#### A The Fourth Paradigm

In relation databases, the fourth paradigm is with the highest paradigm with coexistence of XML functional dependence and XML values dependence [4]. Borrowing the related theories in relation databases, the formal definition of the fourth paradigm (4XNF) is given as below.

Key is an important part in semantic constraints [5]. Key will be set as an important data integrity constraint [6] in the definition of XML mode. The formal definition of key is given below.

**Definition 5(XML key)** Given DTD  $D$ ,  $T$  is a complete tree satisfying  $D$ ,  $p \in \text{Paths}(D)$ . If two arbitrary different nodes  $v_1, v_2$ , ( $v_1$  and  $v_2$  belong to  $N(p)$ ) satisfy  $\text{val}(v_1) \neq \text{val}(v_2)$ ,  $p$  is called as key, and recorded as  $\text{key}(p)$ .

**Definition 6 (XML the fourth paradigm)** Given DTD  $D$ ,  $T$  is a complete tree satisfying  $D$ ,  $\Sigma$  is a XMVD set. If each XMVD in  $\Sigma$ :  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , satisfies any one of below conditions, then DTD  $D$  is the fourth paradigm of XML.

- (1)  $\text{key}(q_i), \text{key}(r_j)$ ;
- (2)  $\text{key}(ph)$  and  $q_i \cap r_j = ph$ ;
- (3)  $\text{key}(ph)$  and  $q_i \cap r_j$  are the strict prefix of  $ph$  ;
- (4)  $q_i \cap r_j = \text{root}$ ;

(5)  $\text{key}(ph)$  and  $r_j \subset q_i$  or  $r_j \subset ph$ ;

where,  $ph = p_1 \cap \dots \cap p_k$ ;  $q_i = q_1 \cap \dots \cap q_m$ ;  $r_j = r_1 \cap \dots \cap r_s$ .

**Lemma 1:** Given DTD  $D$ ,  $\Sigma$  is a XMVD set, and all routes in it are in  $\text{Paths}(D)$ . Assume that  $\Sigma$  causes redundancy, then there exists a multi-value dependence  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , satisfying:

(1)  $v \in N(q_i)$  or  $v \in N(r_j)$ ,  $1 \leq i \leq m, 1 \leq j \leq s$ ;

(2) There exists two different instances  $v_1.v_2. \dots .v_{t-1}.v_t$  ( $v_t = v$ ) and  $w_1.w_2. \dots .w_{t-1}.w_t \in \text{Path}(q_i)$ ,  $1 \leq t$ , which makes  $\text{val}(v_t) = \text{val}(w_t)$ ;

(3) there exists two instances  $z_1, z_2$ ,  $z_1 \in \text{Nodes}(x_{ij}, r_j), z_2 \in \text{Nodes}(y_{ij}, r_j)$ , making  $\text{val}(z_1) \neq \text{val}(z_2)$ ;

(4) there exists two instances  $z_3, z_4$ ,  $z_3 \in \text{Nodes}(x_{ijh}, ph), z_4 \in \text{Nodes}(y_{ijh}, ph)$ , making  $\text{val}(z_3) = \text{val}(z_4), 1 \leq h \leq k$ .

**Proof:** (1) we first prove that if  $\Sigma$  causes redundancy, then there exists a multi-value dependence XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , making it satisfy condition (1) in lemma 1.

If the situation is not established, there are two possibilities. The first one is as below. If  $v \notin N(q_i)$ ,  $v \notin N(r_j)$  and  $v \notin N(ph)$ , because  $v$  does not belong to any node in a route of  $\Sigma$ , when node  $v$  turns to  $v'$ , the generated new tree  $T$  does not violate  $\Sigma$ , and conflicts with  $\Sigma$  to cause redundancy contradiction. The other possibility assumes  $v \in N(ph)$ , in which conditions, if  $\text{lab}(v) \in E_1$ , according definition 1, the changed  $v'$  does not belong to any node in  $T'$ . Same as  $v$  in  $T$ ,  $T'$  satisfies  $\Sigma$  according to definition 3. Therefore,  $\Sigma$  cannot cause redundancy or conflict with  $\Sigma$ . If  $\text{lab}(v) \notin E_1$ , after effective changes, node  $v$  makes  $\text{val}(v')$  not appear in  $T$ . In the new tree  $T$  after the changes, XMVD:  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  still holds, and causes redundancy contradiction with  $\Sigma$ . Thus, condition (1) in lemma 1 is established.

(2) Next we want to prove that if  $\Sigma$  causes redundancy and  $v \in N(q_i)$ , i.e. satisfying condition (1) in lemma 1, there exists a XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , making itself satisfy condition (2) in lemma 1.

Assume this case is not established, then there are two conditions. The first one is: if there is only one route instance in  $\text{Paths}(q_i)$ ,  $\text{Paths}(q_i)$  in the new generated  $T'$  also has only one route instance after node  $v$  experiences an effective change, then XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  is known to be established in  $T'$ , and conflicts with  $\Sigma$ ; if there are more than one routes in  $\text{Paths}(q_i)$  with different tail node values, let  $v_1.v_2. \dots .v_{t-1}.v_t \in \text{Paths}(q_i), v_t = v$ . If in  $\text{Paths}(q_i)$ , there exist route instances  $w_1.w_2. \dots .w_{t-1}.w_t$ , where  $v_t \neq w_t$ , and makes conditions (2) and (3) in definition 7 satisfied, then according to definition 7, it's known that  $v'_1. v'_2. \dots . v'_t-1. v'_t$  are same to  $v_1.v_2. \dots .v_{t-1}.v_t$ . For the same reason,  $w'_1. w'_2. \dots . w'_t-1.w'_t$  are same to  $w_1.w_2. \dots .w_{t-1}.w_t$  according to condition (5) in definition 7, then after effective changes in node  $v$  and

generation of new tree  $T'$ , XMVD: $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  is still established, and conflicts with  $\Sigma$ . The other conditions is: if in  $\text{Paths}(q_i)$  there is no route instances  $w_1.w_2. \dots .w_{t-1}.w_t$ , which satisfy condition (2) and (3) in definition 7, then it's known that after effective changes in node  $v$  and generation of new tree  $T'$ , XMVD: $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  is still established, and conflicts with  $\Sigma$ . Therefore, condition (2) in lemma 1 is established.

(3) Then we prove that if  $\Sigma$  causes redundancy, there exists a multi-value dependence XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , which satisfy condition (3) in lemma 1.

Assume this case is not established. Then there do not exist two nodes  $z_1, z_2$ , where  $z_1 \in \text{Nodes}(x_{ij}, r_j)$ ,  $z_2 \in \text{Nodes}(y_{ij}, r_j)$ , which make  $\text{val}(z_1) \neq \text{val}(z_2)$ . It's known that route  $r_j$  only has one route instance. From definition 7, XMVD: $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  still holds after effective changes in node  $v$  and generation of new tree  $T'$ , XMVD: $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  still holds, and it conflicts with  $\Sigma$ . Thus condition (3) in lemma 1 is established.

(4) Finally, we prove that if  $\Sigma$  causes redundancy, there exists a multi-value XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , which satisfies condition (4) in lemma 1.

Assume this case is not established. Then there do not exist two nodes  $z_3, z_4$ , where  $z_3 \in \text{Nodes}(x_{ijh}, ph)$ ,  $z_4 \in \text{Nodes}(y_{ijh}, ph)$ , which make  $\text{val}(z_3) \neq \text{val}(z_4)$ ,  $1 \leq h \leq k$ . It's known that the tail nodes of the instances in route  $ph$  are different. From definition 7, after effective changes in node  $v$  and generation of new tree  $T'$  XMVD: $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  still holds, and conflicts with  $\Sigma$ . Therefore, condition (4) of lemma 1 is established.

Due to the symmetry property, the cases when  $v \in N(r_j)$  keeps the same conclusions as above. here  $x_{ij}, y_{ij}, x_{ijh}, y_{ijh}, x'_{ij}, y'_{ij}, x'_{ijh}, y'_{ijh}$  holds the same meaning with that in definition 3.

The non-redundancy judgment theorem can be obtained by following definition 6 and lemma 1.

Theorem 2: Given DTD  $D, T$  is a tree satisfying  $D$ ,  $\Sigma$  is a XMVD set,  $T$  satisfies  $\Sigma$ . If  $D$  is with 4XNF mode, then  $\Sigma$  has no redundancy in  $T$ .

Proof: Statements given below prove that the cases of  $v \in N(r_j)$  hold the same results with that in  $v \in N(q_i)$  due to the symmetry property.

(1) Assume that condition (1) in definition 6 holds. That is when  $q_i, r_j$  are both keys,  $\Sigma$  causes redundancy.

According to lemma 1, after node  $v$ 's effective changes into  $v'$ , then there exists  $v_1 \in N(q_i), v_2 \in N(q_i)$ , and  $\text{val}(v_1) \neq \text{val}(v_2)$ . Thus it's known that it conflicts with  $q_i$  being the key, as well as the above assumption.

(2) Assume that condition (2) in definition 6 holds. That is when  $ph$  is a key and  $q_i \cap r_j = ph$ ,  $\Sigma$  causes redundancy.

According to lemma 1, if  $\Sigma$  causes redundancy, that is there exists a complete tree  $T$  satisfying  $D$  and  $\Sigma$ , if there

is a node  $v$  in  $T$ , after effective changes into  $v'$ , a new tree  $T'$  is generated, which leads that XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  does not hold, as well as  $v \in N(q_i)$ . Because XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  is not established in  $T'$ , then  $\text{Paths}(q_i)$  in  $T'$  keep two route instances  $v_1.v_2. \dots .v_{t-1}.v_t$  ( $v_t = v'$ ) and  $w_1.w_2. \dots .w_{t-1}.w_t$ , satisfying condition (1),(2),(3), but not (4) or (5) in definition 7. Because  $ph$  is a key, i.e. the tail node values in instance  $ph$  are different, from condition (3) in definition 7, i.e.  $\text{val}(z_3) \neq \text{val}(z_4)$ , it's known that  $x_{ijh} = y_{ijh}$ ; from condition (4) in definition 7, i.e.  $\text{val}(z'_3) \neq \text{val}(z_3)$ , as well as  $\text{val}(z'_4) \neq \text{val}(z_4)$  in (5), it's known that  $x_{ijh} = x'_{ijh}, y_{ijh} = y'_{ijh}$ . Because  $q_i \cap r_j = ph$ , then we have  $N(q_i \cap r_j) = N(q_i \cap r_j \cap ph) = N(ph)$ , and  $x_{ijh} = y_{ijh} = x'_{ijh} = y'_{ijh} = x_{ij} = y_{ij} = x'_{ij} = y'_{ij}$ . Therefore, if  $T'$  does not satisfy XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , then  $T$  does not satisfy XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , which conflicts with the assumption.

(3) Assume that condition (3) in definition 6 holds. That is when  $ph$  is a key and  $ph$  is a strict prefix,  $\Sigma$  causes redundancy.

According to lemma 1, if  $\Sigma$  causes redundancy, that is there exists a complete tree  $T$  satisfying  $D$  and  $\Sigma$ , if there is a node  $v$  in  $T$ , after effective changes into  $v'$ , a new tree  $T'$  is generated, which leads that XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  does not hold, as well as  $v \in N(q_i)$ . Because XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  is not established in  $T'$ , then  $\text{Paths}(q_i)$  in  $T'$  keeps two route instances  $v_1.v_2. \dots .v_{t-1}.v_t$  ( $v_t = v'$ ) and  $w_1.w_2. \dots .w_{t-1}.w_t$ , satisfying condition (1),(2),(3), but not satisfying (4) or (5) in definition 7. Because  $ph$  is a key, i.e. the node values in instance  $ph$  are different. Further because  $q_i \cap r_j$  is the strict prefix of  $ph$ , then we have  $q_i \cap r_j = q_i \cap r_j \cap ph$ . From condition (2),(3) in definition 7, it's known that  $x_{ij} = y_{ij} = x_{ijh} = y_{ijh}$ . From condition (4),(5) in definition 7, it's known that  $x_{ij} = x'_{ij}, y_{ij} = y'_{ij}, x_{ijh} = x'_{ijh}, y_{ijh} = y'_{ijh}$ . In total,  $x_{ij} = y_{ij} = x'_{ij} = y'_{ij} = x_{ijh} = y_{ijh} = x'_{ijh} = y'_{ijh}$ . Therefore, if  $T'$  does not satisfy XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , then  $T$  does not satisfy XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , which conflicts with the assumption.

Assume that condition (4) in definition 6 holds. That is when  $q_i \cap r_j = \text{root}$ ,  $\Sigma$  causes redundancy.

According to lemma 1, if  $\Sigma$  causes redundancy, which means that there exists a complete tree  $T$  satisfying  $D$  and  $\Sigma$ , if there is a node  $v$  in  $T$ , after effective changes into  $v'$ , a new tree  $T'$  is generated, which leads that XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  does not hold, as well as  $v \in N(q_i)$ . Because XMVD  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m | r_1, \dots, r_s$  is not established in  $T'$ , then  $\text{Paths}(q_i)$  in  $T'$  keeps two route instances  $v_1.v_2. \dots .v_{t-1}.v_t$  ( $v_t = v'$ ) and  $w_1.w_2. \dots .w_{t-1}.w_t$ , satisfying condition (1),(2),(3), but not (4) or (5) in definition 7. Because  $q_i \cap r_j = \text{root}$ , then  $x_{ij} = y_{ij} = x'_{ij} = y'_{ij} = x_{ijh} = y_{ijh} = x'_{ijh} = y'_{ijh}$ , and they are all root nodes. Following condition (3),  $\Sigma$  can be proved to

cause no redundancy, which conflicts with the assumption.

(5) Assume that condition (5) in definition 6 holds, i.e.  $\text{key}(p)$  and  $r_j \subset q_i$  or  $r_j \subset p_h, \Sigma$  causes redundancy. Because  $p_h$  is a key, i.e. the tail node values are different. Then from condition (3) in definition 7, it's known that  $\text{val}(z_3) = \text{val}(z_4)$  and  $x_{ijh} = y_{ij}$ . From  $\text{val}(z'_3) = \text{val}(z_3)$  in (4) and  $\text{val}(z'_3) = \text{val}(z_3)$  in (5), it's known that  $x_{ijh} = y_{ijh} = x'_{ijh} = y'_{ijh}$ . When  $r_j \subset q_i, r_j \cap q_i = r_j$ , and there exists a node  $v$  after effective change to  $v'$ , and generating a new tree  $T'$ , which does not affect  $(r_j)$ . From (4) and (5) in definition 7, it's known that XMVD  $P_1, \dots, P_k \rightarrow Q_1, \dots, Q_m | r_1, \dots, r_s$ . Therefore,  $\Sigma$  does not cause any redundancy, which conflicts with the assumption.

where  $x_{ij}, y_{ij}, x_{ijh}, y_{ijh}, x'_{ij}, y'_{ij}, x'_{ijh}, y'_{ijh}$  have the same meaning to that in definition 3.

### B Design and Algorithm of XML Multi-value Dependence Normalization

Algorithm of XML multi-value dependence normalization is as follows.

Assume that complete XML document tree  $T$  is a XMVD satisfying 3XNF and minimum dependence set  $\Sigma$ .  $P \rightarrow Q | R$  is assumed to be XMVD in  $\Sigma$ . Let  $\{P_1, \dots, P_n\}$  represent instance set of all  $P$  in tree  $T$  satisfying  $P \rightarrow Q | R$ . All responding instance set in  $Q$  is represented as  $\{Q_1, \dots, Q_{qn}\}$ , and all responding instance set in  $Q$  is represented as  $\{R_1, \dots, R_n\}$ .

First of all, for  $P \rightarrow Q | R$ , let  $p$  be the left path public prefix,  $q$  be the public prefix of all routes in  $Q$ ,  $r$  be the public prefix of all routes in  $R$ . For route instance set  $Ph \in \text{Paths}(P)$ , where,  $ph, qh, rh$  are the public prefixes of all route instance sets  $Ph, Qh, Rh$ . If the tail node of  $ph$  is a complex element node, then the longer route between  $qh$  and  $rh$  is chose, which makes the parent node of the route's tail node be  $\text{last}(p)$ ; if the tail node of  $ph$  is a simple element node or attribute node, then the longer route between  $qh$  and  $rh$  is chose, which makes the parent node of the route's tail node be a brother node of  $\text{last}(ph)$ .

For route instance  $P_t$  in  $P$  of different instance  $Ph$ , when  $P_t = Ph$ , we first find out whether  $\text{last}(qh)$  and  $\text{last}(qt)$  have the same value. If not, and the value of  $\text{last}(qt)$  does not equal to the tail node of any route, which have  $\text{pparent}(qh)$  as their parent node, then modify route  $qt$ , and make it be the brother node of  $\text{last}(qh)$ ; when  $\text{last}(qh)$  has same value as  $\text{last}(qt)$ , then delete route  $qt$ . For the same reason, find out whether the value of  $\text{last}(rh)$  equals to that of  $\text{last}(rt)$ , if not, the value of  $\text{last}(rt)$  does not equal to the tail node value of any route with parent node being  $\text{pparent}(rh)$  in  $T'$ , then modify route  $rt$  to be brother node of  $\text{last}(rh)$ ; if  $\text{last}(rh)$  has the same value as  $\text{last}(rt)$ , then delete route  $rt$ . At last, delete route instance set  $P_t$ . Above operations make  $p$  be key route.

After above operations, condition (2) and (3) in definition of 4XNF are satisfied.

Algorithm 2: the algorithm to make DTD  $D$  be changed into  $D'$ , which satisfies 4XFN.

Input: the minimum dependence set  $\Sigma$  of XFD and XMVD, the complete tree  $T$  satisfying  $\Sigma$ .

Output: XML document tree  $T'$  satisfying XMVD paradigm

```

DTD-to-4XNF( $\Sigma, T$ )
Begin
(1)  $T' := T, \Sigma' := \Sigma;$ 
(2) If  $T'$  satisfies 4XNF, then turn to (4);
(3) For each XMVD:  $P \rightarrow Q | R$  do
(4)  $\text{Paths}(P) := \{P_1, \dots, P_n\};$ 
(5)  $P = \{p_1, \dots, p_k\};$ 
(6)  $p = p_1 \cap \dots \cap p_k;$ 
(7)  $\text{Paths}(Q) := \{Q_1, \dots, Q_f\};$ 
(8)  $Q = \{q_1, \dots, q_f\};$ 
(9)  $q = q_1 \cap \dots \cap q_f;$ 
(10)  $\text{Paths}(R) := \{R_1, \dots, R_e\};$ 
(11)  $R = \{r_1, \dots, r_e\};$ 
(12)  $r = r_1 \cap \dots \cap r_e;$ 
(13) For each  $Ph \in \text{Paths}(P)$  do
(14) if any condition in definition of 4XNF is not
satisfied, then
(15) if  $\text{last}(ph) \in E_1$ 
(16) if  $\text{length}(qh) \leq \text{length}(rh)$  then
(17)  $\{\text{vparent}(\text{vparent}(\text{last}(rh))) = \text{last}(ph)\};$ 
(18) Modify  $rh$  in  $T'$  to modified  $rh$ , and modify  $rh$ 
in  $\Sigma'$  to modified  $rh$ ; }
(19) else
(20)  $\{\text{vparent}(\text{vparent}(\text{last}(qh))) = \text{last}(ph)\};$  Modify
 $qh$  in  $T'$  to modified  $qh$ , and modify  $qh$  in  $\Sigma'$  to modified
 $qh$ ; }
(20) else
(21) if  $\text{length}(qh) \leq \text{length}(rh)$  then
(22)  $\{\text{vparent}(\text{vparent}(\text{last}(rh))) = \text{vparent}(\text{last}(ph))\};$ 
Modify  $rh$  in  $T'$  to modified  $rh$ , and modify  $rh$  in  $\Sigma'$  to
modified  $rh$ ; }
(23) else
(24)  $\{\text{vparent}(\text{vparent}(\text{last}(qh))) = \text{vparent}(\text{last}(ph))\};$ 
(25) Modify  $qh$  in  $T'$  to modified  $qh$ , and modify
 $qh$  in  $\Sigma'$  to modified  $qh$ ; }
(26) For route instance  $P_t$  of  $P$  different from  $Ph$  do
(27) If  $(P_t = Ph)$  then
(28) If  $\text{val}(\text{last}(qh)) \neq \text{val}(\text{last}(qt))$  and  $\text{val}(\text{last}(qt))$  is
not equal with the tail node value of any route node with
parent node being  $\text{pparent}(qh)$  then
(29)  $\{x := \text{last}(qt)\};$ 
(30)  $\text{vparent}(x) := \text{vparent}(\text{last}(qh));$ 
(31) Modify  $qt$  in  $\Sigma'$  to modified  $qt$ ; }
(32) else if  $\text{val}(\text{last}(qh)) = \text{val}(\text{last}(qt))$  and
 $\text{vparent}(\text{last}(qh)) \neq \text{vparent}(\text{last}(qt))$  then
(33)  $\text{Paths}(Q) := \text{Paths}(Q) - qt;$ 
(34) if  $\text{val}(\text{last}(rh)) \neq \text{val}(\text{last}(rt))$  and  $\text{val}(\text{last}(rt))$  is
not equal with the tail node value of any route node with
parent node being  $\text{pparent}(rh)$  then
(35)  $\{y := \text{last}(rt)\};$ 
(36)  $\text{vparent}(y) := \text{vparent}(\text{last}(rh));$  Modify  $rt$  in  $\Sigma'$  to
modified  $rt$ ; }
(37) else if  $\text{val}(\text{last}(rh)) = \text{val}(\text{last}(rt))$  and
 $\text{vparent}(\text{last}(rh)) \neq \text{vparent}(\text{last}(rt))$  then
(38)  $\text{Paths}(R) := \text{Paths}(R) - rt;$ 
(39)  $\text{Paths}(P) := \text{Paths}(P) - P_t;$ 
(40) Return( $T'$ );

```

End

Analysis of the algorithm is as follows.

(1) The correctness of the algorithm

Initialization is firstly done to make  $T'$  be  $T$ ,  $\Sigma'$  be  $\Sigma$ . For each multi-value dependence in  $\Sigma$ ,  $P \twoheadrightarrow Q|R$ , let  $p$  be the public prefix of all route in  $P$ ,  $q$  be the public prefix of all route in  $Q$ ,  $r$  the public prefix of all route in  $R$ . Then do the following operations: firstly find out whether this multi-value dependence satisfies the conditions in definition of 4XFN, if not, find whether  $\text{last}(p) \in E1$  is established, if it's established, then judge  $\text{length}(qh) \leq \text{length}(rh)$  is satisfied or not. If it is satisfied, move node in  $T'$  to realize  $\text{vparent}(\text{vparent}(\text{last}(rh))) = \text{last}(p)$ , if it is not satisfied, move node in  $T'$  to realize  $\text{vparent}(\text{vparent}(\text{last}(qh))) = \text{last}(p)$ . The aim of above operations is to make route  $q \cap r$  be the child node of  $p$ . If  $\text{last}(p) \in E1$  is not established, then it's known that this multi-value dependence has only one route in its left, and the tail node is a simple element node or an attribute node, then find out whether  $\text{length}(qh) \leq \text{length}(rh)$  is established or not. If yes, move node in  $T'$  to realize  $\text{vparent}(\text{vparent}(\text{last}(rh))) = \text{vparent}(\text{last}(p))$ ; if not, move node in  $T'$  to realize  $\text{vparent}(\text{vparent}(\text{last}(qh))) = \text{vparent}(\text{last}(p))$ . The aim of above operations is to make route  $q \cap r$  be a brother node of  $p$ .

The next step is to compare  $P_t$  with  $P_h$ . If they are equivalent, then find out whether  $\text{val}(\text{last}(qh)) \neq \text{val}(\text{last}(qt))$  is established or not. If yes, and  $\text{val}(\text{last}(qt))$  does not equal to the tail node value of any route node with parent node being  $\text{pparent}(qh)$  in  $T'$ , then copy tail node of  $qt$  to be a brother node of tail node of  $qh$ . If  $\text{val}(\text{last}(qh)) = \text{val}(\text{last}(qt))$  holds, then delete route  $qt$  in  $T'$ . Moreover, find whether  $\text{val}(\text{last}(rh)) \neq \text{val}(\text{last}(rt))$  is established or not. If yes, and  $\text{val}(\text{last}(rt))$  does not equal to the tail node value of any route node with parent node being  $\text{pparent}(qh)$  in  $T'$ , then copy tail node of  $rt$  to be a brother node of tail node of  $rh$ . If  $\text{val}(\text{last}(rh)) = \text{val}(\text{last}(rt))$  holds, then delete route  $rt$  in  $T'$ . Repeat above operations, and then combine  $q$  and  $r$  in the same  $P$  of  $T'$  to make routes satisfy  $q$  and  $r$  be all strong route. Then condition (1) in definition 6 is satisfied. For any two different routes  $P_t$  and  $P_h$  in  $P$ , if they are equivalent, delete route set  $P_t$  in  $T'$ . All routes that satisfy  $P$  in  $T'$  by this algorithm are strong. Then condition in definition 6 is satisfied.

When the outermost of the cycle is utilized, each XMVD in  $\Sigma'$  satisfies the conditions of the definition 6. Thus  $T'$  obtained by this algorithm meets 4XNF

(2) The capability to be terminated

The algorithm is constructed by three level cycles. The outermost cycle is determined by the number of XMVD, while the two internal cycles are determined by whether complete XML tree is limited. Therefore, the whole algorithm can be terminated because the number of XMVD and XML complete tree are both limited.

(3) The time complexity of the algorithm

The algorithm is constructed by three level cycles. The outermost cycle is determined by the number of XMVD, while the number of the loop of the second cycle is

determined by the maximum rout instance set. Assume that the maximum route instance set contains  $d$  instances, ( $d$  is a multiple of 4). At best, in the first operation of the second cycle, the inetermost cycle operates for  $d-1$  times. In the second operation of the second cycle, the inetermost cycle operates for  $d-5$  times. From this trend, in the last operation of the second cycle, the inetermost cycle operates for 3 times. The second cycle operates for  $d/4$  times at best, it's known that the total times is  $m(d^2+2d)/8$ . At least, for the similar computing process, the total operation times is  $m(d^2-d)/2$ . Therefore, the time complexity is  $O(m \times d^2)$ .

#### IV. EXPERIMENTAL ANALYSIS

##### A Experiment Settings

Experimental environment settings are as follows.

All algorithm realizations in this experiment are as follows.

(1) Hardware ,CPU: Intel 2.20GHz, Memory: 2.00GB, Hard Disk: 250G.

(2) Operation system: Microsoft Windows XP.

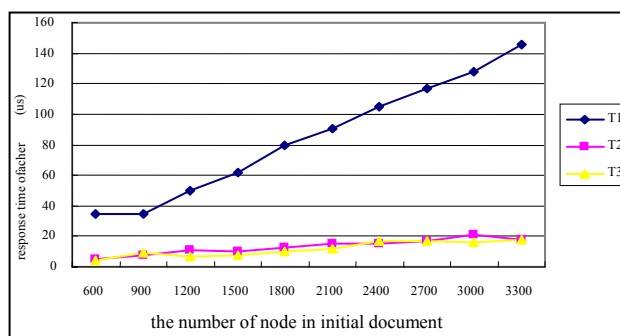
(3) Programming environment: Java, Altova XMLSpy 2010.

Data set settings are as given below . This experiment used combined data set Course and real data set reed to verify the correction of the paradigm of the algorithm.

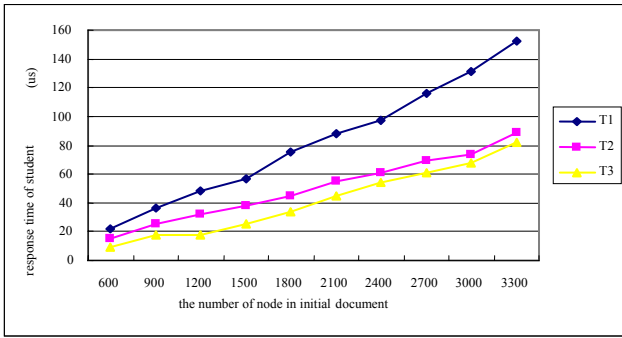
##### B Experiment Result and Performance Analysis

According to the three kinds of DTDs: D1, D2, D3, given by this paper, and the corresponding document trees T1, T2, T3, effectiveness of the normalization algorithm is verified in storage and query efficiency.

Along with the gradual increase of nodes in original document trees T1, T2 are generated after normalization by algorithm 1, and then T3 is generated by algorithm 2. The response time needed for searching a particular node is shown as Figure 3. The experimental results indicate that XML document query response time which satisfies paradigm with higher normalization degree is much lower than that with lower normalization degree.



a) The relationship between response time of teacher and the number of nodes in initial document



b) The relationship between response time of student and the number of nodes in initial document

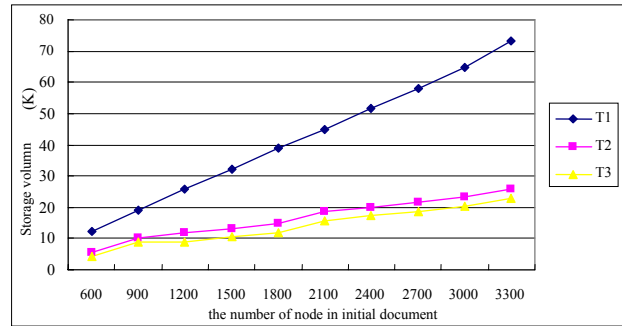
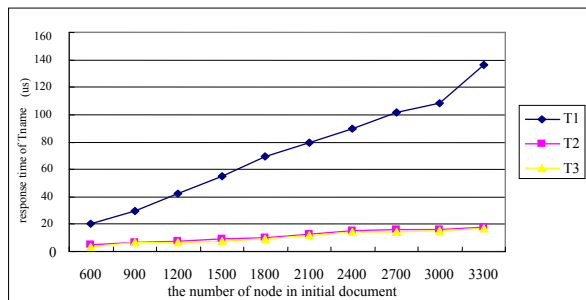
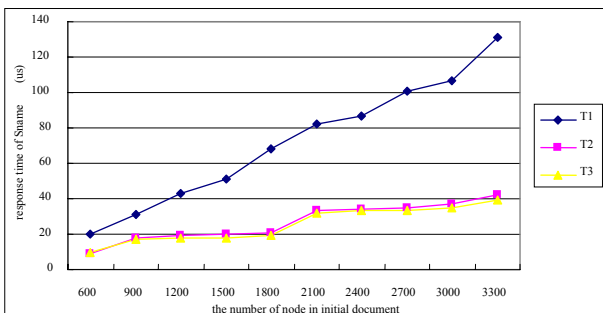


Figure 4 The relationship between storage volume and the number of nodes in initial document



c) The relationship between response time of Tname and the number of nodes in initial document



d) The relationship between response time of Sname and the number of nodes in initial document

Figure 3 The relationship between response time and the number of nodes in initial document

Along with the gradual increase of nodes in original document trees T1, T2 are generated after normalization by algorithm 1, and then T3 is generated by algorithm 2. The storage needed for the three document trees is as shown in Figure 4. The experimental results indicate that XML document storage which satisfies paradigm with higher normalization degree is much lower than that with lower normalization degree.

The real data base uses the XML data base of course information of Reed College. T2 is generated based on original document tree T1 after normalization by algorithm 1, and then T3 is generated by algorithm 2. The query response time of the three document trees is shown as Figure 5. The experimental results indicate that XML document query response time which satisfies paradigm with higher normalization degree is much lower than that with lower normalization degree.

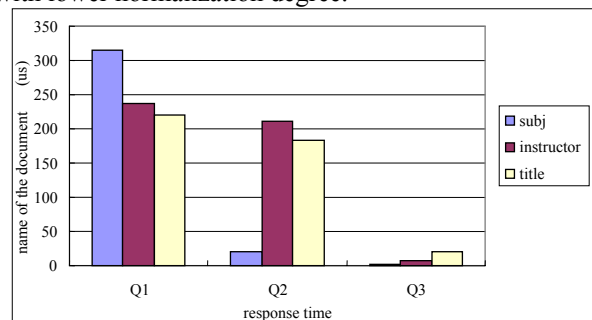


Figure 5 The relationship between response time and the number of node in initial document

T2 is generated based on T1 after normalization by algorithm 1, and then T3 is generated from T2 by algorithm 2. The storage needed of the three documents is as shown in Figure 6. The experimental results show that XML document storage which satisfies paradigm with higher normalization degree is much lower than that with lower normalization degree.

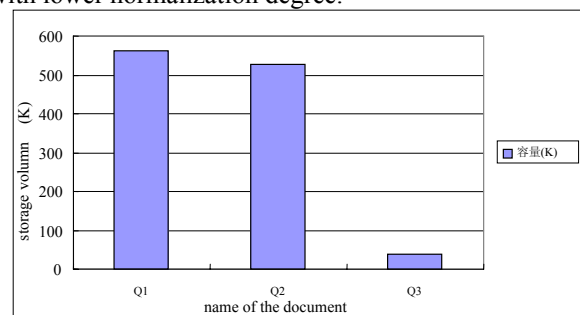


Figure 6 The relationship between storage volume and the number of node in initial document

V CONCLUSIONS



This paper describes the redundancy and key in XML, and presents the formal definition of 3XNF and 4XNF. Normalization rules and non-redundancy judgement theorem are then proposed. Furthermore, normalization algorithm based on XML document is also given which is followed by corresponding analysis and proof with respect to not only the capability to be terminated, but also the correctness, as well as the computation complexity. Finally, it has been shown that the effectiveness of the algorithm can be demonstrated by the experiments.

#### ACKNOWLEDGEMENT

This work has been financially supported by the Project of Major Reform to the Network Engineering Programme in Hebei Normal University of Science and Technology, which was assigned by Hebei Education Department in 2012.

This work has been also sponsored by Qinhuangdao science and technology research (201101A038\201101A185\2012021A133)

#### REFERENCES

- [1] James J. Lu, Sonali Renjen. Normalizing XML Schemas through Relational Database[C]//Proceedings of the 43rd Annual Southeast Regional Conference, Kennesaw, Georgia, Alabama, USA, 2005: 220-221.
- [2] Cong Yu, H.V.Jagadish. XML schema refinement through redundancy detection and normalization[J]. the VLDB Journal, 2008, 17(2): 203-223.
- [3] Flavio Ferrarotti, Sven Hartmann, Henning Köhler, et al. The Boyce-Codd-Heath Normal Form for SQL[C]//Proceedings of 18th International Workshop on Logic, Language, Information and Computation, Philadelphia, PA, USA, 2011:110-122.
- [4] Flavio Ferrarotti, Sven Hartmann, Henning Köhler, et al. Foundations for a Fourth Normal Form over SQL-Like Databases[C]//Conceptual Modelling and Its Theoretical Foundations. Essays Dedicated to Bernhard Thalheim on the Occasion of His 60th Birthday. London, USA, 2012:85-100.
- [5] Md. Sumon Shahriar, Jixue Liu. Towards the Preservation of Keys in XML Data Transformation for Integration[C]//Proceedings of the 14th International Conference on Management of Data, IIT Bombay, Mumbai, India, 2008: 116-126.
- [6] Flavio Ferrarotti, Sven Hartmann, Schastian Link, et al. Performance Analysis of Algorithms to Reason about XML Keys[C]//Proceedings of 23rd International Conference of Database and Expert Systems Applications, Vienna, Austria, 2012:101-115.
- [7] Guoliang Li, Beng Chin Ooi, Jianhua Feng, Jianyong Wang, Lizhu Zhou. EASE: An Effective 3-in-1 Keyword Search Method for Unstructured, Semi-structured and Structured Data[C]// Proceedings of the ACM SIGMOD international conference on Management of data ,SIGMOD 2008, Vancouver, BC, Canada, 2008: 903-914.
- [8] Richard Cyganiak, David Wood. RDF 1.1 Concepts and Abstract Syntax[EB/OL]. (2011-08-30)[2012-10-15].<http://www.w3.org/TR/2011/W-D-rdf11-concepts-20110830>.
- [9] Flavio Ferrarotti, Sven Hartmann, Schastian Link. A Precious Class of Cardinality Constraints for Flexible XML Data Processing[C]//Proceedings of 30th International Conference on Conceptual Modeling, Brussels, Belgium. 2011:175-188.
- [10] M. W. Vincent, J. Liu, M. Mohania. The implication problem for 'closest node' functional dependencies in complete XML documents[J]. Journal of Computer and System Sciences, 2012, 78(4): 1045-1098.
- [11] Jixue Liu, Jiuyong Li, Chengfei Liu, Yongfeng Chen. Discover Dependencies from Data—A Review[J]. IEEE Transactions on Knowledge and Data Engineering, 2012, 24(2):251-264.
- [12] Sven Hartmann, Sebastian Link, Thu Trinh. Solving the Implication Problem for XML Functional Dependencies with Properties[C]//Proceedings of the 28th international Workshop on Logic, Language, Information and Computation, Brasilia, Brazil, 2010:161-175

**Lijun Cao**, born in 1971, has a master degree. She is an associate professor in the college of mathematics and information science of Hebei Normal University of Science and Technology. her main research directions include data mining, grid technology.

**Xiyin Liu**, born in 1970, has a master degree. He is an associate professor in the college of mechanical and electrical engineering of Hebei Normal University of Science and Technology. His main research directions include data mining, grid technology.