INFORMATION TO USERS

This was produced from a copy of a document sent to us for microfilming. While the most advanced technological means to photograph and reproduce this document have been used, the quality is heavily dependent upon the quality of the material submitted.

The following explanation of techniques is provided to help you understand markings or notations which may appear on this reproduction.

- 1. The sign or "target" for pages apparently lacking from the document photographed is "Missing Page(s)". If it was possible to obtain the missing page(s) or section, they are spliced into the film along with adjacent pages. This may have necessitated cutting through an image and duplicating adjacent pages to assure you of complete continuity.
- 2. When an image on the film is obliterated with a round black mark it is an indication that the film inspector noticed either blurred copy because of movement during exposure, or duplicate copy. Unless we meant to delete copyrighted materials that should not have been filmed, you will find a good image of the page in the adjacent frame.
- 3. When a map, drawing or chart, etc., is part of the material being photographed the photographer has followed a definite method in "sectioning" the material. It is customary to begin filming at the upper left hand corner of a large sheet and to continue from left to right in equal sections with small overlaps. If necessary, sectioning is continued again-beginning below the first row and continuing on until complete.
- 4. For any illustrations that cannot be reproduced satisfactorily by xerography, photographic prints can be purchased at additional cost and tipped into your xerographic copy. Requests can be made to our Dissertations Customer Services Department.
- 5. Some pages in any document may have indistinct print. In all cases we have filmed the best available copy.

University Microfilms International

300 N. ZEEB ROAD, ANN ARBOR, MI 48106 18 BEDFORD ROW, LONDON WC1R 4EJ, ENGLAND POURNAGHSHBAND, HASSAN

A NEW APPROACH FOR CONSTRUCTING A RELATIONAL SCHEMA FROM A SET OF DATA DEPENDENCIES

The University of Oklahoma

PH.D. 1980

.1

.

University Microfilms International 300 N. Zeeb Road, Ann Arbor, MI 48106

.THE UNIVERSITY OF OKLAHOMA

· GRADUATE COLLEGE

A NEW APPROACH FOR CONSTRUCTING A

RELATIONAL SCHEMA FROM A SET OF DATA DEPENDENCIES

A DISSERTATION

SUBMITTED TO THE GRADUATE FACULTY

in partial fulfillment of the requirements for the

degree of

DOCTOR OF PHILOSOPHY

BY

HASSAN POURNAGHSHBAND

Norman, Oklahoma

A NEW APPROACH FOR CONSTRUCTING A RELATIONAL SCHEMA FROM A SET OF DATA DEPENDENCIES

APPROVED BY in Icno DISSERTATION COMMITTEE

To my lovely mother, Haji Khanoom.

.

ACKNOWLEDGMENTS

I wish to express my sincere appreciation and deep gratitude to Dr. John C. Thompson, for his keen interest, capable guidance, and invaluable encouragement during my graduate program and particularly throughout the planning and completion of this study.

I am also sincerely thankful to the members of my dissertation committee, Dr. S. Lakshmivarahan, Dr. John T. Minor, Dr. Albert B. Schwarzkopf, and Dr. Bill K. Walker, for their helpful and valuable suggestions.

It is also my pleasure to thank Ms. Betty Sudduth for her expert typing of this manuscript.

My special gratitude to my family, Haji Khanoom, Nahid and Ahmad, Roohangiz and Hossein, Krista and Reza, Pouri and Mehdi, Hayedeh and Javad, and Massoud, for their patience and continuing encouragement throughout my academic career.

Finally, my especial thanks go to someone who has made my life worthwhile.

iii

ABSTRACT

It is the purpose of this work to investigate the Third and Fourth normal form relational schemas.

The relational model and the notion of relational data bases were first presented by Codd. Codd introduced First, Second, the Third normal forms, and presented an approach to the 3NF decomposition. He defined a normalized relation as one for which each of the underlying domains contains atomic values only, so that every value in the relation is in turn atomic.

In addition to Codd's decomposition approach, another approach to the logical design of relational data bases was the synthetic approach. Among several attempts to this approach, Bernstein's algorithm is efficient and has been proved to be correct.

A new normal form (4NF) was proposed by Fagin in 1977. Fagin presented a decomposition approach to the 4NF relations. His decomposition procedure leads to a family of 4NF relations which is not necessarily "optimal".

The purpose of this work is to present a new approach to constructing 4NF relations from functional dependencies

iv

and multivalued dependencies. The objectives of the procedure are two-fold: To make the task as algorithmic as possible, and to produce an "optimal" 4NF family.

TABLE OF CONTENTS

																		Page
ABSTRA	CT .		• •	• •	•	•	•	•	•	•	•	•	•	•	•	•	•	iv
LIST O	F ILLU	USTRATI	ONS	• •	•	•	•	•	•	•	•	•	•	•	•	•	•	ix
Chapte	r																	
I.	INTRO	ODUCTIO	N .	••	•	•	•	•	•	•	•	•	•	•	•	•	•	1
	$1.1 \\ 1.2$	Histor The Re	y . lati	ona]	Mo	ode	1	•	D	ef	in	i.t	ic	ons	•	•	•	1
		and No	tati	ons	•	•	•	•	•	•	•	•	•	•	•	•	•	6
	1.3	Operat	ions	on	Re.	lat	:ic	ns		•	•	•	•	•	•	٠	•	9
	1.4	Normal	izat	ion	•		•	•	•	•	٠	•	•	•	•	٠	•	12
	1.5	Data D	epen	denc	cies	5	•	•	•	•	•	•	٠	٠	•	•	•	15
	1.6	Normal	For	ms .	•	٠	•	•	•	•	٠	•	•	٠	٠	•	٠	18
11.	THE A MULT	ALGEBRA IVALUED	OF DEP	FUNC ENDI	CTIC ENCI	DNA CES	AL S	DE •	PE •	ND	EN •	•	ES.	5 I •	•	•	•	22
	2.1 2.2	Introd The In	ucti fere	on . nce	Ru	Les	• s f	or	D	at	• a	•	•	•	•	•	•	22
		Depend	enci	es.	•	•	•	•	•	•	•	•	•	3	٠	•	•	22
		2.2.1 2.2.2	Arm Fun Inf	stro ctic erer	ong' onal nce	's L I Ru	Ax Dep 11e	io en s	ma de fo	ti nc	za ie Mu	uti es ult	or iv	n c val	of • .ue	•d	•	23
		2.2.3	Dep Inf	ende	enci nce	Les Ru	; ile	•	fo	• or	Mi	• .xe	d	•	•	•	•	26
			Dep	ende	enc:	ies	3	•	•	•	•	•	•	•	•	•	•	<u>3</u> 0
	2.3 2.4 2.5	The Cl The Pr Chapte	osur inci r Su	es c ples mmai	of I s o: cy a	Dep E S anc	oer Sch I F	ide iem Rem	nc a ar	ie De ƙs	s si	.gr	• •	•	• •	•	•	32 36 3 <u>8</u>
III.	CONS DEPEI	TRUCTIN NDENCIE	GRE S.	LATI	ION	sc •	CHE •	ME •	s •	FR	•	1 I •	PAC •	Α'. •	•	•	•	-39
	3.1 3.2	Introd Genera	ucti 1 de	on . scri	ipt:	ior	• 1S	of	• t	he	A		orc	• ac	• h	•	•	· 39 40

.

Chapter

.

	3.3	A Simpl Method	Le Depe	ende	ncy-	·to-	Rela •••	tic •••	on •	•	•	•	•	41
		3.3.1 3.3.2	Proof Proof	of of	Repr Norn	ese nali	ntat zati	ion	۰ ا	•	•	•	•	43 44
	3.4	The Dep lst Imp	pendenc	cy-t ent	0-Re	elat •	ion	Met •••	hoc.	1, •	•	•	•	46
		3.4.1	Proof	of	Repr	ese	ntat	ior	ı.	•	•	•	•	46
	3.5	The Dep 2nd Imp	pendenc	∵y-t ent	0-Re	elat •	ion	Met	hod.	1, •	•	•	•	51
		3.5.1	Proof	of	Repi	ese	ntat	ior	ı.	•	•	•	•	52
	3.6	The Dep 3rd Imp	pendenc	∷y-t ent	0-Re	elat	ion	Met	hoo.	1, •	•	•	•	57
		3.6.1 3.6.2	Proof Proof	of of	Repr Norn	ese Mali	ntat zati	cior Lon	1.	•	•	•	•	58 62
	3.7	The Dep 4th Imp	pender. proveme	cy-t ent	0-Re	elat	ion	Met	hoo.	đ,	•	•	•	65
		3.7.1 3.7.2	Proof Proof	of of	Repi Nori	ese nali	ntat zati	cior Lon	۱. •	•	•	•	•	70 71
	3.8	The Dep Final 1	pendeno Improve	ey-t	o-Re	elat	ion	Met	hoo•	а, •	•	•	•	75
	3.9	Chapter		ary 	and	Rell		•	•	•	•	•	•	04
IV.	CONST	PRUCTING	3 CLOSI	JRE	11 6	AND	CLOS	ORE	; 1.	LT	•	•	•	85
	4.1 4.2	Introdu Descrip	action otion o	of t	he C	los	ure	II	Alq	goj	it	hn	• 1.	85 85
		4.2.1 4.2.2 4.2.3 4.2.4	Proof Proof Output Speed	of of : An Ana	Tern Corn alys lysi	nina rect sis ls	i ior ness	1. 5.	•		• • •	•	• • •	89 89 92 95
	4.3	Descrip Algori	otion other that the second se	oft	he (los •	ure	II] • •	•	•	•	•	•	95
		4.3.1 4.3.2 4.3.3 4.3.4	Proof Proof Output Speed	of of : An Ana	Cori Tern alys lysi	rect nina sis Ls	ness tior	5 . 1 .	•	•	• •	•	• • •	96 101 101 102

Page

Chapter

,

	4. 4 4. 5	Descrij Algori Chapter	ption of thms r Summary	the Do and I	ectec Remar	tior ks	1 • •	••	•	•	•	•	103 103
۷.	ANALY RELAT	YSIS OF FION .	RELATION	S BY 1	DEPEN	IDENC	CY-7	-07	•	•	•	•	106
	5.1 5.2	Introdu Propert	uction . ties of R Algorith	elations	ons C	Const	ruc	tec		•	•	•	106 106
		5.2.1 5.2.2 5.2.3	Represen Separatio Minimal	tation on Pri Redund	n Pri incip dancy	ncip ple 7 Pri	ole Inci		•	• •	• •	•	107 113 117
			5.2.3.1	Desc: Deco	ripti nposi	on c tior	of t NAl	he goi	Op cit	oti :hn	.mz າ	ıl	129
	5.3 5.4 5.5	Bernste Fagin's Chapter	ein's Syn s Decompo r Summary	thetion sition and 1	c App n App Remar	oroac oroac :ks	ch. ch.	•	•	• •	• •	• •	131 134 137
VI.	SUMMA	ARY AND	CONCLUSI	on .	•••	••	• •	•	•	•	•	•	138
	6.1 6.2 6.3	Summary Analys: Sugges	y is of the tions for	Appro Furtl	bach ner W	 Iork	• • • •	• •	•	•	•	•	138 139 141
BIBLIO	GRAPHY	Y	● ල ù ♦	• • •	•••	• •	• •	•••	•	•	•	•	143
APPEND	EX .		• • • •	• • •	• •			•	•	•	•	•	146

Page

.

LIST OF ILLUSTRATIONS

Figure	•	Pa	ge
1-1	An example of a relation	•	2
1 - 2 a	An example of relational schema (consisting of three relational schemes)	•	8
l-2b	A snapshot of the schema of part "a"	•	8
1-3	The projection operation	•	
1-4	The join operation	•	11
1 - 5a	An example of a relation with undesirable properties	•	13
1-5b	A different normal form of the schema of part (a), without undesirable properties .	•	13
1 - 6a	A third normal form schema with undesirable properties	•	16
1-6b	The fourth normal form of the schema of part "a", without undesirable properties .	•	16
1-7	Transitive Dependency	•	20
3-1	Algorithm 3-1	•	42
3-2	Algorithm 3-2	•	47
3-3	Algorithm 3-3	•	53
3-4	Algorithm 3-4	•	59
3-5	Derivation of $A \rightarrow C$ from $A, D \rightarrow B, C$ and $A \rightarrow D$.	•	64
3-6	Algorithm 3-5	•	66
3-7	Algorithm 3-6	•	76
4-1	Algorithm 4-1	•	87
4-2	Algorithm 4-2	•	97

.

4-3	Algorithm 4-3	•	•	•	٠	•	•	٠	.104
51a	An instance of the schema R .	•	•	•	•	•	•	•	. 109
5-1b	An instance of the schema S .	•	•	•	•	•	•	•	. 109
5-2	An instance of the schema S'.	•	•	•	•	•	•	•	. 111
5-3	Algorithm 5-1	•	•	•	•	•	•	•	.130
5-4	Algorithm 5-2	•	•	•	•	•	•	•	. 132

~

CHAPTER I

INTRODUCTION

1.1 History

One of the most significant contributions to data management technology in recent years has been the development of the relational point of view of a data base. In this section we will give a very informal discussion of the relational approach and its problems and will trace its historical developments.[†]

The concept of the relational model of data was first proposed in 1970 by Codd [10]. Conceptually, a relation can be viewed as a table, such as the one in Figure 1-1. Note that the table resembles the traditional file, with rows corresponding to records of the file and columns corresponding to fields of the records [13]. This tabular representation of data makes the computer accessible not only to the professional users, but also to the casual users with little or even no programming background. The conceptual simplicity of the relational model together with its mathematical elegance have attracted a large number of followers. This

[†] Formal definitions and technical discussions are given later.

HOUSING

TENANT	COMPLEX	TYPE	RENT	APT#	MANAGER
Jack	Nieman	Fur.	100	17 D	Smith
Mary	Parkview	Unfur.	110	212H	Brown
Phil	Kraettli	Fur.	200	305F	Rogers
Mark	Kraettli	Fur.	200	302A	Rogers
Mike	Kraettli	Fur.	200	208C	Rogers
John	Const.	Unfur.	175	126H	Rose
Tom	Const.	Unfur.	175	113F	Rose

Figure 1-1

dissertation deals mainly with the relational model, and specifically it considers the problem of "logical schema design", i.e. how relations should be designed to take advantage of the relational model's strengths.

Shortly after Codd introduced the relational model for data bases, he observed that certain relationships (functional dependencies) contained in a relation can cause data manipulation anomalies (i.e., adding new information, modifying existing information, or deleting old ones might be difficult and sometimes impossible). Therefore, he proposed a hierarchical set of constraints on relations. These constraints restrict relations to certain canonical forms, called Normal forms by Codd. The most desirable of these is Third normal form (3NF)[11]. Subsequently several researchers proposed various algorithms for constructing 3NF relations [8,15,29].

All of the above works are based on the axiomatization of functional dependencies developed by Armstrong [3]. Wang and Wedekind proposed an algorithm to synthesize the relational schema from a set of functional relationships [29]. However, their algorithm is not correct, in the sense that it could generate two relations with keys that are functionally equivalent.

Another approach was made by Bernstein [8]. His work was based on the approach given by Delobel and Casey [15]. Although the relational schema produced by Bernstein's

algorithm is in 3NF, it may contain properties that can cause data manipulation anomalies. This problem, which was discovered by Fagin [18] and independently by Zaniolo [31], stems from the fact that there may exist a special type of relationship between attributes which is not functional (a so-called multivalued dependency). Considering this new concept of relationship, Fagin defined a new normal form which he called fourth normal form (4NF). Then he proposed a decomposition approach for constructing 4NF relations [18]. His process begins by forming a single relation consisting of all attributes of the data base. Then this relation is broken down into two subrelations in a more desirable form (i.e., causing less data anomalies). This process continues for each subrelation not in 4NF, until all of them are in The family of 4NF relations constructed in this way is 4NF. not necessarily "optimal", and may contain redundant infor-In addition, for a large data base, the original mation. relation produced by Fagin's approach can be extremely large (it consists of all attributes of the data base), and thus makes the problem (i.e., breaking the relations) costly and time consuming.

The main objective of the new approach given in this thesis is also to construct 4NF relations. In this approach (in contrast to Fagin's method), the original and final relations are approximately of the same size. This is in fact the primary advantage of the approach. Briefly, the process

begins by constructing one relation for each relationship that exists between any pair of attributes (or sets of attributes). Then these relations are broken down to subrelations until all of them are in 4NF. The relational schema constructed by this approach contains the theoretical minimal number of attributes and thus optimizes the use of computer memory.

The organization of the thesis is as follows: In this chapter, basic definitions and notations are given. Also a detailed discussion about normalization and normal forms are presented.

In <u>Chapter II</u>, the theoretical bases for designing a relational data model in general, and for our new approach in particular, are discussed.

The algorithm for constructing a relational schema in Fourth Normal Form is presented in <u>Chapter III</u>. The chapter starts with a simple algorithm which produces a schema not necessarily in 4NF. Then this simple algorithm is modified to eliminate undesirable properties of the schema as much as possible. This modification process is continued until we come up with an algorithm (which in this dissertation is called Fourth Normal Form Algorithm, or briefly, FNF) which produces a schema in 4NF. Theoretical issues and all proofs regarding the above algorithms are also described in Chapter III.

In <u>Chapter IV</u>, two algorithms are introduced for finding Closure II and Closure III (two sets that play significant

roles in our approach) for a given set of data dependencies. The feasibility of these algorithms and proofs for their correctness are also discussed in this chapter.

In <u>Chapter V</u>, properties of relations constructed by the FNF algorithm are examined and they are compared with relations constructed by other approaches. Also in this chapter a detailed discussion about "optimal" schemas is given.

Finally, in Chapter VI, an investigation of the practical and theoretical consequences of our approach is given. Also, several directions for further research are pointed out at this time.

The proofs of those problems related to this research that have been dealt with in previous work by others are given in the <u>appendix</u>.

1.2 The Relational Model -- Definitions and Notations

Conceptually, a <u>relation</u> can be viewed as a table in which each row corresponds to a distinct <u>entity</u> (or <u>tuple</u>) and each column to a distinct attribute. There exists a set of possible associated values, called the <u>domain</u> of attribute for each attribute.

Formally, a <u>relation</u> can be defined as follows: Given a collection of sets D_1, D_2, \ldots, D_n (not necessarily distinct), a relation of R, defined over D_1, D_2, \ldots, D_n , is a subset of the cartesian product $D_1 \times D_2 \times \ldots \times D_n$. That is, R is a set of n-elements each of the form (d_1, d_2, \ldots, d_n) where $d_i \in D_i$, for $1 \le i \le n$. Each element of the R is called a <u>tuple</u> of R. R is said to be a relation of degree n. An attribute is a

name assigned to a domain of a relation. While the domains of a relation need not be distinct, the attribute names assigned to them must be unique within the relation.

Suppose A_1, A_2, \ldots, A_n are the names of the domains D_1, D_2, \ldots, D_n of a relation R, then we use notation 1.2.1 for R.

$$R(A_1, A_2, \dots, A_n)$$
 (1.2.1)

The attribute set of R is defined as

$$U = \{A_1, A_2, \dots, A_n\}$$
 (1.2.2)

We will also use (1.2.3) to designate R on the set of attributes U.

R(U) (1.2.3)

The structure of a relation is sometimes called the <u>intention (scheme)</u> and the contents of a relation is referred to as the <u>extension</u>. The contents of a relation may vary from time to time. That is, tuples may be modified, deleted from a relation, or/and inserted in a relation. The contents of a relation in a particular time is called its <u>snapshot</u> (or <u>instance</u>). Figure 1-2 indicates a schema (i.e., a collection of schemes) consisting of three relation schemes, $R_1(A,B)$, $R_2(C,D,E)$ and $R_3(B,D,F)$; and a snapshot of the schema.

Let R(U) be a relation on the set of attributes U, and let W be a subset of U, then W is a <u>candidate key</u> of R if it can uniquely identify the tuples of R and no proper subset of W has this property. A relation may

A Relational Schema and its Snapshot

a) An example of relational schema (consisting of three relational schemes)

$$R_{1}(A,B) = \{(a_{1},b_{2}), (d_{2}b_{4}), (d_{3},b_{1})\}$$

$$R_{2}(C,D,E) = \{(c_{1},d_{1},e_{2}), (c_{2},d_{3},e_{1}), (c_{3},d_{2},e_{1}), (a_{1},d_{3},e_{3})\}$$

$$R_{3}(B,D,F) = \{(b_{1},d_{2},f_{3}), (b_{2},d_{2},f_{2}), (b_{3},d_{1},f_{1})\}$$

b) A snapshot of the schema of part "a"

d₂

a_l

f2

fl

^b2

b₃

R ₁ (A,	B)		$R_2(C)$	D _i	E)
al	^b 2		cl	al	^e 2
^a 2	b ₄		°2	a3	el
^a 3	^b l		°3	₫ ₂	e _l
			°4	a3	e ₃
R ₃ (B,	D,	F)			
b ₁	a ₂	f ₃			

Figure 1-2

have more than one candidate key. An attribute is said to be <u>prime</u> if it appears in any candidate key of the relation, and it is called a <u>non-prime</u> attribute otherwise.

1.3 Operations on Relations

Codd has defined [10] a number of operations on relations. In this dissertation two operations are of particular interest: Projection and Join.

Let R be a relation defined on the set of attributes $U = A_1, A_2, \dots, A_n$. For any $W = A_1, A_2, \dots, A_m$, that is $W \subseteq U$, the projection of R on W is defined as:

$$\Pi \neq R(W) = (a_1, a_2, \dots, a_m) \quad (a_1, a_2, \dots, a_n) \in R \quad (1.3.1)$$

In other words, we can think of the projection of R onto W as the operation that takes the relation (instance) represented by R, then deletes all columns except those labeled by attributes in W, and finally identifies common tuples (See Figure 1-).

The join operator which is in some senses an inverse to the projection operator, in fact connects attributes of different relations together. The join (natural join) of a relation R(X,Y) with a relation P(Y,Z)on Y, is defined as:

 $R * P = \{(x,y,z) | (x,y \in R \text{ and } (y,z) \in P\}.$ (1.3.2) Figure 1- indicates an example of join operation.

The projection operation

Suppose an instance of relation R(A,B,C) is

R(A,	В,	С,	D,	E)
al	b ₃	°2	d ₁	e3
al	b ₂	°3	d _l	e ₄
a ₂	b4	°ı	đ ₃	e ₄
a ₃	bl	°2	đ ₂	e ₂
a3	bl	°3	ď2	el

Then the projection of R on A, B and D is

R'(A,	В,	D)
al	b ₃	ďl
al	b ₂	ďl
a2	b ₄	ď3
a ₃	b ₁	d ₂

The join operation

If instances of relations R(A,B,C) and P(C,D) are

R(A,	в,	C)	P(C,	D)
al	b ₃	°2	c _l	ďl
al	^b 2	°3	°2	ď3
^a 2	b ₄	cl	c ₄	ď2
^a 3	bl	°2		
^a 3	b ₁	c ₄		

then the natural join of R with P (i.e., R \star P) is

RP(A,	в,	с,	D)
al	b ₃	°2	a3
a ₂	b4	°1	a _l
a ₃	b _l	°2	d3
a ₃	b ₁	c ₄	^d 2

.

1.4 Normalization

The notion of normalization in relational data base was first presented by Codd [11]. Codd observed that certain relations have structural properties that are undesirable for describing data bases. These under sirabilities stem from the fact that some attributes in a relation are related to each other in a certain way. For example, in relation TCM of Figure 1-5a, there is a relationship between COMPLEX and MANAGER. That is, given the complex, we can determine its manager. Note that in relation TCM the association between a complex and its manager is repeated for each TENANT. This repetition causes data manipulation anomalies. These anomalies can be categorized as follows:

> (1) Update anomaly. Suppose the association between a complex and its manager is changed (i.e., the complex is no longer managed by the old manager, but a new one). Then we need to update all tuples and change the values of complex and manager for all tenants of the complex. On the other hand, if we update only one tuple, it will not be adequate for maintaining a consistant schema. This is known as update anomaly.

a) An example of a relation with undesirable properties.
 (The association between a complex and its manager is repeated for each tenant.)

TCM (TENANT,	COMPLEX,	MANAGER)
Jack	Nieman	Smith
Mary	Parkview	Brown
Phil	Kraettli	Rogers
Mark	Kraettli	Rogers
Mike	Kraettli	Rogers
John	Const.	Rose
Tom	Const.	Rose

b) A different normal form of the schema of part "a", without undesirable properties.

TC (TENANT,	COMPLEX)	CM (COMPLEX,	MANAGER)
Jack	Nieman	Nieman	Smith
Mary	Parkview	Parkview	Brown
Phil	Kraettli	Kraettli	Rogers
Mark	Kraettli	Const.	Rose
Mike	Kraettli		
John	Const.		
Tom	Const.		

Figure 1-5

- (2) Insertion anomaly. If a new TENANT moves into a COMPLEX, and he/she happens to be the first tenant of the complex, then an association between the complex and manager will also be needed to insert the new tuple. This is called insertion anomaly.
- (3) Deletion anomaly. If a TENANT moves out from a complex, and he/she happens to be the last tenant of the complex (i.e., no more tenants are residing in the complex), then by removing his tuple from the relation, the association between the complex and its manager also disappears from the relation. This is known as deletion anomaly.

Data manipulation anomalies discussed above will disappear if we convert the schema into a "better" normal form. (See Figure 1-5b).

For a relational schema, in fact, there are usually different sets of relation schemes that can represent all of the information needed. Thus, to avoid data manipulation anomalies attempts have been made to introduce schemas with no undesirable structural properties for describing data bases. These considerations led Codd to define a series of three normal forms [11], First Normal Form, Second Normal Form, and Third Normal Form.

Later in 1977, Fagin [18] discovered that even by putting a schema into Third Normal Form, not all of the anomaly problems necessarily disappear. This led him to propose a new normal form, called Fourth Normal Form. An example of a Third Normal Form schema with undesirable properties is illustrated in Figure 1-6a. Figure 1-6b shows the Fourth Normal Form of the schema.

Data manipulation anomalies appear in this third normal form schema, because of the undesirable relationship that exists between attributes EMPLOYEE and CHILD in relation EDC.

A formal description of all normal forms is given in section 1.6.

1.5 Data Dependencies

The problems associated with constructing Fourth Normal Form relational schemas are tied to the fact that some attributes determine the values of other attributes. This can be formalized as the comcept of Functional Dependency (FD) and Multivalued Dependency (MVD).

The concept of functional dependency is defined as follows: if X and Y are distinct collections of attributes of some relation scheme R, then Y is said to be functionally dependent on X (or X functionally determines Y) if, at every point in time, each X value has no more than one Y a) A Third Normal Form schema with undesirable properties. (The association between an employee and his/her department is repeated for each child.)

EDC (EMPLOYEE,	DEPARTMENT,	CHILD)
Johnson	Accounting	Henry
Johnson	Accounting	Bonnie
Johnson	Accounting	Ralph
Stewart	Personnel	Tim
Stewart	Personnel	Martha
Jones	Marketing	Joe

b) The Fourth Normal Form of the schema of part "a", without undesirable properties.

ED (EMPLOYEE,	DEPARTMENT)	EC (EMPLOYEE,	CHILD)
Jo hnson	Accounting	Jo hnson	Henry
Stewart	Personnel	Johnson	Bonnie
Jones	Marketing	Johnson	Ralph
		Stewart	Tim
	·	Stewart	Martha

Joe

Jones

$$R.X \longrightarrow R.Y \tag{1.4.1}$$

and if no confusion exists, then R can be omitted.

$$X \longrightarrow Y$$
 (1.4.2)

The concept of Multivalued Dependencies has been proposed by Fagin in [18] as follows:

Let $R(X_1, \ldots, X_m, Y_1, \ldots, Y_n, Z_1, \ldots, Z_r)$ be a relation with m + n + r attribute names. For notation convenience, we write X for X_1, \ldots, X_m ; Y and Z are defined analogously. Whenever we write, say, R(X, Y, Z), we assume automatically X, Y, and Z are pairwise disjoint as above. If x_1, \ldots, x_m are entries that appear under columns X_1, \ldots, X_m , then we write x for (x_1, \ldots, x_m) ; y and z are defined analogously. Define Y_{χ_Z} to be $\{y: (x, y, z) \in R\}$. Of course Y_{χ_Z} is nonempty if and only if x and z appear together in a tuple of R. Now, using the following notation to denote that X multidetermines Y in R

$$R.X \longrightarrow R.Y$$
 (1.4.3)

and when no confusion exists

$$X \longrightarrow Y$$
 (1.4.4)

the multivalued dependency $X \longrightarrow Y$ is said to hold for R(X,Y,Z) if Y_{XZ} depends only on x; that is, if $Y_{XZ} = Y_{XZ}$, for each x,z,z' such that Y_{XZ} and Y_{XZ} , are nonempty. As we can see the validity of MVD $X \longrightarrow Y$ depends not only on the values of attributes X and Y, but also on the value of Z, the complement of XY. This context dependency of MVDs makes the problem of schema design much more complicated than if only FDs were involved.

As an example of Multivalued dependencies see relation EDC of Figure 1-6, for which multivalued dependency $EMPLOYEE \longrightarrow CHILD$ holds.

1.6 Normal Forms

The concepts of functional dependence and multivalued dependence play significant roles in the theory which governs the decomposition of relations into subrelations in normal forms.

To show how a certain undesirable dependency creates problems discussed earlier, we will discuss the concepts of partial dependencies (and fully dependencies) and transitive dependencies mentioned by Codd [10,11]. We will also discuss the Fagin's [18] notion of nontrivial multivalued dependency.

We say that, Y is <u>fully dependent</u> on X in relation R, if

- X and Y are two distinct subcollections of attributes of relation R,
- (2) $R.X \longrightarrow R.Y$, and
- (3) Y is not functionally dependent on any proper subset of X.

If condition (3) is not satisfied, then we say, Y is partially dependent on X in relation R.

Partial dependencies can create data manipulation anomalies and thus, have to be removed from the schema. This led Codd to further normalize the first normal form relations to get the second normal form [11].

A relation R is said to be in <u>First Normal Form</u> (1NF), if and only if all underlying domains contain atomic values only [13].

A relation R is said to be in Second Normal Form if:

- (1) it is in first normal form, and
- (2) every non-prime attribute of R is fully dependent on each candidate key of R.

To define third normal form relations, we need to know the concept of "transitive dependencies".

Given a relation R, further suppose that X, Y, and Z are three distinct collections of attributes of R, and if the following conditions are true

- (1) $R.X \longrightarrow R.Y$
- (2) R.Y \rightarrow R.X
- $(3) \quad R.Y \longrightarrow R.Z$

then it follows that

$$R.X \longrightarrow R.Z$$
 and

 $R.Z \longrightarrow R.X$

Here, Z is said to be <u>transitively dependent</u> on X under the relation R. This concept is depicted in Figure 1-7.



 \rightarrow given FDs \implies implied FDs

Figure 1-7

A relation is said to be in Third Normal Form (3NF) if,

- (1) it is in second normal form, and
- (2) every non-prime attribute of R is nontransitively dependent on each candidate key of R.

The concept of "trivial multivalued dependencies" proposed by Fagin [18], is needed in describing the fourth normal form relations.

Given relation R(U) where $U = \{X,Y\}$, then the multivalued dependencies $X \longrightarrow \emptyset$ and $X \longrightarrow Y$ necessarily hold for R. These are called <u>Trivial Multivalued</u> Dependencies (TMD).

A relation R is said to be in Fourth normal Form (4NF), if whenever a nontrivial multivalued dependency $X \longrightarrow Y$ holds for R, then so does the functional dependency $X \longrightarrow A$ for every attribute A of R. An example of a schema in 4NF is given in Figure 1-6.

CHAPTER II

THE ALGEBRA OF FUNCTIONAL DEPENDENCIES AND MULTIVALUED DEPENDENCIES

2.1 Introduction

The theoretical bases for designing a relational data model in general, and for our new design approach in particular, are discussed in this chapter. The inference rules for functional dependencies and multivalued dependencies are given in section 2.2, and it is proved that these rules are complete in the sense that a data dependency g is a consequence of a set of dependencies G if and only if g can be derived from G by a sequence of applications of the rules.

In section 2.3 different closures of dependencies are defined, and it is shown that two of these closures are of particular importance for this work.

Finally, the three principles that should be considered in designing a relational schema are presented in section 2.4. These principles are formally discussed in Chapter V.

2.2 <u>The Inference Rules for Data Dependencies</u> When the problem of schema design is of concern,

it is important to know whether a dependency is implied by some other dependencies or not. Formally, this means that a dependency g is a consequence of a set of dependencies G if for all relation schemes R, g holds in R if all dependencies of G hold in R. In previous work, researchers have proposed and investigated the complete sets of inference rules for functional and multivalued dependencies [3,6,7,18, et al.]. Detailed discussions for these inference rules are given in the following subsections:

2.2.1 Armstrong's Axiomatization of Functional Dependencies

Axiomatization of functional dependencies was studied by Armstrong [3]. In [3] he has presented a set of axioms governing sets of functional dependencies. It is proved [3,7,19] that this set of axioms is complete for the family of FDs in the sense that, for a family of FDs, for each set F of FDs from the family, the FDs that are implied by F are exactly those dependencies that can be inferred from it using these inference rules. In other words, we say a set of axioms is complete for a family of FDs if and only if, for each set F of FDs over a set of attributes U, if F^{\dagger} is the set of FDs that follow logically from F, then every relation over U that satisfies F, also satisfies the functional dependencies in F⁺. For any formal system, the completeness of the
set of inference rules is an important concept for the system. For the family of functional dependencies, as it is mentioned in [7], only if a complete set of axioms is used, the data base designer can be assured that he has complete knowledge of all FDs that hold in the data base. This is in fact the basic reason that the completeness of Armstrong's axioms has been an important basis for research in this area (including the present research).

The complete set of axioms for the family of functional dependencies is given below. In the rules, X, Y, Z and W are arbitrary subsets of U, where U is the set of all attributes. We write X for the set X, and XY for the union of two arbitrary sets.

FD Rules:

FDl (Reflexivity): If $Y \subseteq X$ then $X \to Y$. FD2 (Augmentation): If $Z \subseteq W$ and $X \to Y$, then $XW \to YZ$. FD3 (Transitivity): If $X \to Y$ and $Y \to Z$, then $X \to Z$.

Other Useful Rules:

FD4 (Pseudo-Transitivity): If $X \rightarrow Y$ and $YW \rightarrow Z$, then $XW \rightarrow Z$.

FD5 (Union): If $X \rightarrow Y$ and $X \rightarrow Z$, then $X \rightarrow YZ$. FD6 (Decomposition): If $X \rightarrow YZ$, then $X \rightarrow Y$ and $X \rightarrow Z$. If α and β are attributes of a relation R, then by applying the Axiom FDl to $X = {\alpha, \beta}$ we get $\alpha\beta \rightarrow \alpha\beta$, $\alpha\beta \rightarrow \alpha, \ \alpha\beta \rightarrow \beta, \ \alpha \rightarrow \alpha$, and $\beta \rightarrow \beta$.

Axiom FD2 means that, knowing f: $X \rightarrow Y$, we can construct another functional dependency, say g: $X, \delta \rightarrow Y$, where the attributes appearing on the left side of g consists of X plus some other extraneous attributes δ , whose values have no effect on the value of Y selected by g.

For Axiom FD3, assume one is given the dependencies $f: X \rightarrow Y$ and $g: Y \rightarrow Z$. The axiom claims that there is a dependency h: $X \rightarrow Z$. Symbolically, h(X,Z) is defined to be g(f(X),Z).

Notice that, in the above set of axioms, the Axioms FD1-FD3 are sufficient, and the other three axioms, that is FD4-FD6 are implied by the first three axioms. As an example, Axiom FD4 can be derived from the Axioms FD1-FD3 as follows: As our assumption we have $f_1: X \rightarrow Y$ and $f_2: YW \rightarrow Z$. Now, from f_1 and Axiom FD2 we get $f_3: XW \rightarrow YW$. By applying Axiom FD3 to f_2 and f_3 we can derive an FD $XW \rightarrow Z$, completing the claim. Similarly, it is easy to show that the other two axioms, FD5 and FD6, can also be derived from the first three axioms.

Let F be a set of functional dependencies over a set of attributes U. Then the <u>closure</u> of F, denoted by F^+ , is defined to be the set of all functional dependencies that can be obtained by successive application of Axioms FDL,

FD2, and FD3 on the set F, over the set of attributes U. Notice that by Armstrong's theory, if F is a given set of FDs for a relation R, then each FD in F^+ also exsits in R.

A functional dependency $f_i \in F$ is said to be <u>redun-</u> <u>dant</u> in F if $F^+ = (F - f_i)^+$. Also, H is a <u>nonredundant</u> <u>covering</u> of a given set of functional dependencies F, if $F^+ = H^+$ and H contains no redundant FDs. As we will see later, the concepts of closures and covering are important in constructing the relational schemas from dependencies. Letting F^+ be the general closure of a given set of dependencies, we will present in Section 2.3 two special closures of a set of functional and multivalued dependencies. These two closures are of particular importance in constructing a relational schema from a set of data dependencies using the new design approach discussed in this research.

2.2.2 Inference Rules for Multivalued Dependencies

A set of inference rules for multivalued dependencies is presented in [7], and it is proved that the given set is complete for the family of MVDs. Similar to FDs, we say that a set of inference rules is complete for the family of MVDs, if for each set M of MVDs from the family, the MVDs that are implied by M are exactly those dependencies that can be inferred from it using these inference rules. The complete set of rules for multivalued dependencies is given below. In the rules, X, Y, Z, and W are arbitrary sets of attributes. We use XY for the union of two arbitrary sets of attributes X and Y. <u>MVD Rules</u>:

MVDO (Complementation): If U = XYZ and $Y \cap Z \subseteq X$, then $X \longrightarrow Y$ if and only if $X \longrightarrow Z$. MVD1 (Reflexivity): If $Y \subseteq X$ then, $X \longrightarrow Y$. MVD2 (Augmentation): If $Z \subseteq W$ and $X \longrightarrow Y$ then, $XW \longrightarrow YZ$. MVD3 (Transitivity): If $X \longrightarrow Y$ and $Y \longrightarrow Z$ then, $X \longrightarrow Z - Y$.

other useful rules:

MVD4 (Pseudo-Transitivity): If $X \rightarrow Y$ and $YW \rightarrow Z$ then, $XW \rightarrow Z - YW$.

MVD5 (Union): If $X \rightarrow Y$ and $X \rightarrow Z$ then, $X \rightarrow YZ$.

MVD6 (Decomposition): If $X \rightarrow Y$ and $X \rightarrow Z$ then, $X \rightarrow Y \cap Z$, $X \rightarrow Y - Z$, and $X \rightarrow Z - Y$.

It is interesting to note that for each FD rule there exists one MVD rule corresponding to it. But notice, that there is no FD rule corresponding to the complementation rule of MVDs. In fact, as it was mentioned in Chapter I, multivalued dependencies are sensitive to context, while functional dependencies are context independent. In other words, to know if the MVD $X \rightarrow Y$ holds in a relation R(U), where U = XYZ, we also need to know the values of the attributes in the set U - XY, that is, Z. As is discussed in [18], an MVD $X \rightarrow Y$ may not hold in one relation, but in one of its projections. This is not true for the case of FDs, because if an FD $X \rightarrow Y$ holds in a relation R(U), it also holds in any other relation R(U') as long as $X \bigcup Y \subset U'$.

As for FD rules, the inference rules MVD0-MVD3 are sufficient for multivalued dependencies, and the other useful rules indicated above (i.e., MVD4-MVD6), can be derived by rules MVD0-MVD3. This claim can be formalized as follows [24]:

<u>Proposition 2.1</u> The Decomposition Rule (i.e., MVD6) can be derived from MVD1-MVD3.

<u>Proof</u>: Suppose we have $m_1: X \to Y$ and $m_2: X \to Z$, then we need to prove that (i) $X \to Y \cap Z$ and (ii) $X \to Z - Y$. From MVD m_1 and Axiom MVD2 we get $m_3: X \to XY$, and from m_2 and Axiom MVD2 we get $m_4: XY \to Z$. Now by applying Axiom MVD3 to m_3 and m_4 we can derive $m_5: X \to Z - XY$. We can also get $m_6: X(Z - XY) \to Z$ and $m_7: X \to X(Z - XY)$ by applying Axiom MVD2 to m_2 and m_5 respectively. From MVDs m_6 and m_7 , and Axiom MVD3 we can get $m_8:$ $X \to Z - X(Z - XY)$. Now, it is easy to get $m_9:$ $X \to Y \cap Z - X$ from m_8 (by boolean manipulations) which in turn can result to $m_{10}: X(X \cap Y \cap Z) \to Y \cap Z$, completing the claim for part (i). For part (ii), we may apply Axiom MVD2 to m_5 to get m_{12} : $Z((X-Y) \cap Z) \longrightarrow (X-XY)$ $((X-Y) \cap Z)$ which in turn can lead us to m_{13} : $X \longrightarrow Z - Y$ (boolean manipulations), thus completing the proof. #

As we saw in Proposition 2.1, the Decomposition Rule can be derived without using the Complementation Rule MVDO. In fact, in [7] it was concluded that neither Rule MVD5 nor Rule MVD6 can be derived without using the complementation rule. This claim has been rejected in [24] for the Rule MVD6, because as we saw in Proposition 2.1 the Rule MVD6 is obtainable without using the Rule MVD0. But notice that, the Union Rule can be derived only if the Complementation Rule is also involved in the derivation process [7,24]. This claim is formalized in the following proposition [24].

<u>Proposition 2.2</u> The Union Rule (i.e., MVD5) cannot be derived from Rules MVD1-MVD3.

Proof: The proof is given in the Appendix.

The significance of the complementation rule is more formally stated by Mendelzon [24] and others [3,18]. Indeed, it has been proved that the complementation rule does not follow from MVD1-MVD6.

In [24], Mendelzon has also proved that there are only two minimal complete subsets of the MVD Rules, the sets {MVD0, MVD1, MVD3} and {MVD0, MVD1, MVD4}. These are important results that can be of particular interest in

designing a relational schema. We will return to this problem in Chapter III.

Recall that a complete set of inference rules for functional dependencies was given in section 2.2.1, and a complete set of inference rules for multivalued dependencies was given in this seciton. Now, it is important to notice that, although FD rules are sufficient when there are only FDs, and MVD rules are sufficient when when there are only MVDs, the combination of FD rules and MVD rules is not sufficient for the case that both . kind of dependencies, functional and multivalued, are involved. In other words, when we have a set of functional and multivalued dependencies, not all of the obtainable dependencies from this set can be inferred by the applications of FD rules and MVD rules. Indeed, there are two additional rules that can be applied only when both kind of dependencies exists in the set. These rules are given in the following subsection.

2.2.3 Inference Rules for Mixed Dependencies

In subsections 2.2.1 and 2.2.2 we dealt with the inference rules for FDs and MVDs respectively. In fact, in each subsection we concentrated on one type of dependency. That is, we wanted to know what additional FDs (or MVDs) can be implied by a set of FDs (or MVDs). Now, suppose G = F [] M is a set of data dependencies, where F is a set of FDs and M is a set of MVDs. Also, suppose

that F' is a set of FDs inferred from F by applying FD rules, and M' is a set of MVDs inferred from M by the applications of MVD rules. Now the question is, does the set $G' = F' \bigcup M'$ contain all dependencies implied by G? If not, what are the additional rules that can be used to deduce them?

In [7] three rules are introduced for mixed dependencies. They are given below. In the rules X, Y, Z, and Z' arbitrary sets.

FD-MVD Rules;

FD-MyD1: If $X \rightarrow Y$ then $X \rightarrow Y$.

FD-MVD2: If $X \longrightarrow Z$ and $Y \longrightarrow Z'$ where $Z' \subseteq Z$ and $Y \cap Z = \emptyset$, then $X \longrightarrow Z'$.

additional useful rule:

FD-MVD3: If $X \longrightarrow Y$ and $XY \longrightarrow Z$, then $X \longrightarrow Z - Y$.

Note that, the first two rules (FD-MVDl and FD-MVD2) combined with the FD rules and the MVD rules are sufficient for mixed dependencies, and the Rule FD-MVD3 can be derived from them. The complete proofs concerning these rules are given in the appendix. Briefly, the Rule FD-MVDl simply follows from the definitions of functional and multivalued dependencies, and it means that an FD is also an MVD. Indeed, as it was mentioned in Chapter I, functional dependencies are special cases of multivalued dependencies (notice that the converse of this concept is not true). As we saw, the set of rules introduced in subsections 2.2.1, 2.2.2, and 2.2.3 are complete for the family of functional and multivalued dependencies. The proof of this completeness which originates from [7] is given in the appendix.

2.3 The Closures of Dependencies

If $G = F \bigcup M$ is a set of data dependencies, where F is a set of FDs and M is a set of MVDs, then the <u>closure</u> of G denoted by G^+ is the set of FDs and MVDs that can be deduced from F $\bigcup M$ by repeated applications of FD rules, MVD rules, and mixed rules.

Recall from the previous section that the set of rules {FD1,FD2,FD3,MVD0,MVD1,MVD2,MVD3,FD-MVD1,FD-MVD2} is sufficient, and thus the axioms of this set are those that are crucial in designing a relational schema. In fact, as we saw before, the other axioms are implied by those given in the above set, and need not be explicitly considered in the design process. In addition, we claim that Axioms FDl and MVDl are not needed either for this This is simply because FD X \rightarrow X (or MVD X \rightarrow X) work. means, having a set of attributes X, we can determine (or multidetermine) a set of attributes X, and clearly, having a set X at our disposal we do not have to determine it. We also claim that Axioms FD2 and MVD2 are not needed either. In fact, if a set of attributes X

determines (or multidetermines) Y, then our method realizes that any set which contains X also determines (or multidetermines) Y.

Axioms MYD0 and FD-MVD1 need not be considered (explicitly) either, because as we will see in Chapter III, these rules are implicitly considered in our schema design.

As will be discussed in more detail in Chapter III, the above concepts lead to two special closures that are of particular interest for our new design approach. We formally define these closures as follows: <u>Definition 2.1</u> If $G = F \bigcup M$ is a set of mixed dependencies (FDs and MVDs), the <u>closure II</u> of G, denoted by G^{II} , is defined to be the set of all dependencies that can be obtained by successive application of axioms FD4 and MVD4.

An efficient closure II Algorithm which computes the closure II of a given set of FDs and MVDs is given in Chapter IV. The feasibility of constructing the closure II of a set for the real world applications is also discussed in that chapter.

<u>Definition 2.2</u> Let $G = F \bigcup M$ be a set of mixed dependencies (FDs and MVDs), then the <u>closure III</u> of G, denoted by G^{III} , is defined to be the set of all dependencies that can be obtained by successive application of axioms FD4, MVD4, and FD-MVD2.

In chapter IV, an efficient algorithm which computes

the closure III of a given set of functional and multivalued dependencies is given, and its feasibility is discussed in detail.

In all the published work of relational data base design [4,8, et al.], the computation of the closure of a set of dependencies has been avoided. This is mainly because the closure of a set of dependencies can be extremely large even if the original set is small. In particular it is the Reflexivity and Augmentation Rules, which produce a large number of dependencies from a given set. On the other hand, the closure II and closure III of a set of dependencies will be fairly small, and therefore we need much less time to compute them than if we were to compute the general closure.

Here, an example is given to clarify these concepts, and a formal discussion with all necessary proofs is given in Chapter IV.

Example 2.1: Given G, the following set of FDs and MVDs:

```
 \begin{array}{ll} \mathbf{f}_1 \colon & \mathbf{A} \longrightarrow \mathbf{B} \\ \mathbf{f}_2 \colon & \mathbf{B} \longrightarrow \mathbf{D} \\ \mathbf{m}_1 \colon & \mathbf{C} \longrightarrow \mathbf{D} \end{array}
```

we can construct G^+ , the closure of G as follows: Applying reflexivity rules we get

$$f_{3}: A \to A$$
$$f_{4}: B \to B$$
$$f_{5}: C \to C$$

$$f_{6}: D \rightarrow D$$

$$m_{2}: \Lambda \rightarrow A$$

$$m_{3}: B \rightarrow B$$

$$m_{4}: C \rightarrow C$$

$$m_{5}: D \rightarrow D$$

applying augmentation rules we get

 $f_{7}: A, C \longrightarrow B$ $f_{8}: A, D \longrightarrow B$ $f_{9}: B, A \longrightarrow D$ $f_{10}: B, C \longrightarrow D$ $m_{6}: C, A \longrightarrow D$ $m_{7}: C, B \longrightarrow D$

applying transitivity rules we get

 $f_{11}: A \rightarrow D$

applying complementation rule we get

 $m_8: C \rightarrow A, B$

applying Axiom MyD1 we get

 $m_{g}: A \longrightarrow B$ $m_{10}: B \longrightarrow D$

and finally applying Axiom MVD2 we get

$$f_{12}: C \rightarrow D$$

Now, we need to get the union of this new set. That is, we have to remove those dependencies that are repeated in the set. This is actually what makes the problem costly and time consuming.

The same process should be applied to this new set to get more dependencies (if any). We continue this process until no more newer set is derivable.

As is indicated in the above Example, constructing the closure of a set of dependencies is very costly and time consuming (a formal treatment to this problem is given in Chapter IV. On the other hand, if we want to construct the closure III (and/or closure II), then we need only to apply the transitivity rules and FD-MVD2. If we do so, we get $G^{III} = \{A \rightarrow D, C \rightarrow D\}$, and $G^{II} =$ $\{A \rightarrow D\}$. By comparing G^+ with G^{III} and G^{II} , for the same set of dependencies, we can observe how much faster an algorithm can construct the closure II (and/or closure III) of a set of dependencies, than if it was to construct the general closure of the set.

2.4 The Principles of Schema Design

When constructing a relational schema from a set of data dependencies, there are principles that have to be considered. In other words, since the intention is to construct schemas with no structural properties that are undesirable for describing data bases (these undesirable properties were discussed in Section 1.4), we should consider these principles when using the dependencies for our schema design. Therefore, we say, a schema is a desirable

one if it satisfies these principles.

The three principles of the schema design that are of interest for this work are proposed in [6]. These three principles are representation, separation, and nonredundancy.

Roughly speaking, representation means that all information of interest should somehow be represented in the schema. On the other hand, if any information is lost during the transformation of dependencies (i.e., the process of schema design), then the constructed schema violates the representation principle, and thus it can no longer be a desirable schema.

A formal discussion is given in Chapter V. We will give our own definitions of these principles, and will prove that the schema constructed by our approach satisfies them.

By separation, it is meant that in designing a relational schema, attempts should be made to represent the basic information separately. This is, in fact, the motivation behind the normalization process [13]. In Chapter V, we will formally discuss this issue, and will prove that the schema constructed by our algorithm is as normalized as possible, that is, it is in fourth normal form.

The nonredundancy principle (sometimes called minimal redundancy) requires that the representation be

nonredundant, that is, while the schema must represent all of the information of interest, but it should not contain any redundant information. Specifically speaking, dependencies that can be derived from other explicitly represented dependencies, need not be explicitly represented. In addition, attributes should be repeated in as few relations as possible. In Chapter V, heuristics are employed for this issue, and it is formally shown how a design approach can lead to an "optimal" schema.

2.5 Chapter Summary and Remarks

The theoretical bases for designing a relational schema have been discussed in details. A set of axioms (FD rules) for the family of functional dependencies has been given, and it is shown that these axioms are complete for this family. Also, a complete set of rules (MVD rules) has been presented in this chapter for the family of MVDs. It has been stated that the combination of FD rules and MVD rules is not sufficient for the family of functional and multivalued dependencies, and thus additional rules (FD-MVD rules) have been given to complete the set of rules for FDs and MVDs.

Furthermore, the three design principles, i.e., representation, separation, and nonredundancy have been introduced.

CHAPTER III

CONSTRUCTING RELATION SCHEMES FROM DATA DEPENDENCIES

3.1 Introduction

<u>.</u>

A design approach for constructing a relational schema from a set of functional and multivalued dependencies is examined in this chapter. Since the intention is to propose an approach satisfying the three design principles discussed in Chapter II, then it is shown in this chapter that the schema constructed by this method:

- (i) "represents" <u>all</u> information of interest in the schema,
 - (ii) does not represent any "redundant" information and
 - (iii) <u>is</u> as "normalized" as possible, that is, it is in fourth normal form.

The proofs concerning these principles, and the proofs of the other related problems are also given in this chapter (except the proofs for the final algorithm -- these proofs will be given in Chapter V).

It is shown, that although the constructed schema has no structural properties that are undesirable for describing

our data base, but it is not necessarily "optimal". This issue will be examined in details in Chapter V. We will propose the heuristics that can be employed for designing an algorithm which can construct an optimal schema.

3.2. General Descriptions of the Approach

We start our work with designing a simple algorithm and show several undesirable properties of the schema constructed by this method. That is, we will prove and also show by examples that some of the design principles are not satisfied, and thus may cause data manipulation anomalies. Then we will modify the algorithm to eliminate those undesirable properties as much as possible. We will continue this process of modification until we come up with an algorithm which produces a "good" schema. This "good" schema will not necessarily be "optimal", and thus we will further modify the algorithm (in Chapter V) to construct a schema which is both, "good" and "optimal".

Following each algorithm one or more examples are given to clarify the concepts considered (and/or missed) by the algorithm.

Also, we will put into account some basic semantics properties that attributes and relationships among them (FDs and MVDs) apparently have in the real world.

As a real world example we will consider a "university housing" data base, which is used as a test case throughout

the iterative refinement of the algorithm.

3.3. A Simple Dependency-to-Relation Method

The problem of the schema design is indeed: given a set of dependencies (FDs and MVDs designated by the data base administrator), how to construct a relational schema. For this method, we present an algorithm which gets as input a set of dependencies, and generates as output a set of relation schemes with designated keys (Figure 3-1).

The algorithm simply constructs one relation scheme for each explicitly given dependency (and not for dependencies that can be inferred from the given set), and designates the keys as follows: (i) if the corresponding dependency is an FD, then the key of the relation scheme would be the attributes appearing on the left side of the dependency, and (ii) if the corresponding dependency is an MVD, then the key of the scheme would be all of the attributes appearing in the dependency.

As we will see later, the schema constructed by this algorithm has some structural properties that are not desirable for describing the data base. As we mentioned before, these undesirable properties will be eliminated by further modification of the algorithm.

Note that, for all of the algorithm presented in this chapter, we assume without loss of generality that there are no two FDs f_i and f_j with identical left sides. In fact,

Algorithm 3-1

A Simple Algorithm to Construct a Schema From a

Set of Dependencies

INPUT: A set F of m FDs; and a set M of n MVDs.

OUTPUT: A set R of m + n relation schemes (in 1NF or better). <<construct relation schemes from FDs>>

<u>do</u> for each $f_i \in F$;

construct a relation scheme R_i consisting of all attributes appearing in f_i ; The set of attributes that appears on the left side of f_i is the key of the relation scheme;

end;

$$\begin{array}{c}
 m \\
 R' = \bigcup & R_{i}; \\
 i = 1
 \end{array}$$

<< construct relation schemes from MVDs>>

```
do for each m_i \in M;
```

construct a relation scheme R_{i+m} consisting of all attributes appearing in m_i ; The set of all attributes

that appears in m_i is the key of the relation scheme; end;

$$R'' = \bigcup_{\substack{i=m+1}}^{m+n} R_i;$$

<<pre><<put all schemes together>>

R = R'UR'';

end algorithm;

42

Figure 3-1.

if there are any, we can simply combine them, using Axiom FD5.

3.3.1. Proof of representation

A schema embodies a set of dependencies (FDs and MVDs) if each dependency is embodied in at least one relation of the schema. An FD, g: $X \longrightarrow Y$ is embodied in a relation $R(A_1, A_2, \dots, A_n)$ if [8]:

- (i) $X \subseteq \{A_1, A_2, ..., A_n\},$
- (ii) $Y \subseteq \{A_1, A_2, \dots, A_n\}$, and
- (iii) X ε DOM(X) and y ε DOM(Y) implies that g(x) = y if and only if the tuple {x,y} is in the projection of R on X and Y.

As the reader will notice, it is easy to generalize the idea to MVDs.

<u>Theorem 3.1</u>. Let R be a relational schema constructed from a set of dependencies G = FUM, using the algorithm 3-1, then R represents the same information as G (i.e., embodies F and M).

<u>Proof</u>: The algorithm constructs one relation (scheme) for each $f_i \in F$ and one relation for each $m_i \in M$. Thus we need only to prove that each of those relations embodies its corresponding dependency. Assume that the relation scheme $R_i(X,Y)$ is constructed by dependency $g_i: X \longrightarrow Y$ (or $X \longrightarrow Y$), then obviously $X \subseteq \{X,Y\}, Y \subseteq \{X,Y\}$. And since the projection of $R_i(X,Y)$ on X and Y is the R_i itself, thus x ε DOM(X) and y ε DOM(Y) <u>implies</u> that $g_i(x) = y$ if and only if the tuple {x,y} is in R_i , completing the proof. #

3.3.2. Proof of Normalization

The schema constructed by the simple algorithm, clearly is not in the desirable normal form (i.e., 4NF), but as it is formally stated below, it is in some degree of normal forms.

<u>Theorem 3.2</u>. Let R be a relational schema constructed from a set of dependencies, using the Algorithm 3-1, then R is in first normal form.

Proof: The proof follows directly from the definition of first normal form. #

As we saw in the previous subsections, our simple algorithm satisfies the representation principle, and can only guarantee a schema of 1NF. We will show by some examples that the nonredundancy principle and 4NF schema cannot be guaranteed by this algorithm.

Example 3.1: Suppose the following data dependencies are given by the data base administrator:

 $f_1: A \longrightarrow B, C$ $m_1: D, E \longrightarrow F$ $m_2: D \longrightarrow E, F$

then, the algorithm constructs one relation scheme for each FD or MVD as follows: (underscored attributes are keys) $R_1(\underline{A}, B, C)$ $R_{2} (\underline{D}, \underline{E}, \underline{F})$ $R_{3} (\underline{D}, \underline{E}, \underline{F})$

and after removing relation ${\rm R}^{}_3$ (or ${\rm R}^{}_2)$ from the schema, we get

- $R_1(\underline{A}, B, C)$
- $R_2(\underline{D},\underline{E},\underline{F})$

Although the algorithm works well for the cases like what we had in the above example, but it may not be appropriate for other situations. The following example clarifies this concept.

Example 3.2: In our UNIVERSITY HOUSING data base, if each COMPLEX is managed by only one MANAGER, and if each MANAGER can manage a set of COMPLEXes, then we may have the following dependencies and their corresponding relation schemes:

f1:	$COMPLEX \longrightarrow$	MANAGER	Rl	(COMPLEX,	MANAGER)
-----	---------------------------	---------	----	-----------	----------

 $m_1: MANAGER \longrightarrow COMPLEX R_2 (MANAGER, COMPLEX)$

As we can see, R_1 and R_2 bear the same information (in the sense of attributes), and R_1 represents not only what it must represent, but also what is represented by R_2 (in the sense that in R_1 not <u>all</u> of the attributes make the key of the relation). Thus, R_2 is redundant and must be removed from the schema.

We can eliminate the above problem, and all problems of this type (but not necessarily the violation of nonredundancy principle), simply by searching for relation schemes with the same set of attributes and "stronger" relationship/s among them (recall that an FD is an stronger dependency than an MVD). Formal discussion for this issue is given in Chapter V. Here we will solve the problem by modifying our simple algorithm. This modification is considered in PART-2 of the following algorithm (Figure 3-2). Notice that, this part searches only the set of schemes R' constructed from FDs. It does not consider the set R" constructed from MVDs because, it is known from PART-1 of the algorithm, that there are not any pair of schemes in the set R" that have the same set of attributes.

3.4. The Dependency-to-Relation Method, 1st Improvement

As it was discussed in the previous subsection, the simple algorithm (Algorithm 3-1) constructs a schema which may have undesirable properties. The algorithm presented in this section (Algorithm 3-2) eliminates some of these undesirable properties, but not all of them. Thus the algorithm will be further modified in the subsequent subsections.

3.4.1. Proof of Representation

<u>Theorem 3.3</u>. Let R be a relational schema constructed from a set of dependencies G = FUM, using the Algorithm 3-2, then R represents the same information as G.

<u>Proof</u>: The algorithm first constructs one relation scheme for each dependency, and then removes a relation scheme R_i

Algorithm 3-2

An Algorithm to Construct a Schema From a Set of

Dependencies; 1st Improvement

INPUT: A set F of m FDs; and a set M of n MVDs.

OUTPUT: A set R of relation schemes (in 1NF or better).

Part-1: <<this part constructs one relation scheme for

each dependency>>

<< construct relation schemes from FDs>>

<u>do</u> for each f_i ε F;

construct a relation scheme R_i consisting of all attributes appearing in f_i ; The set of attributes that appears on the left side of f_i is the key of the relation scheme;

end;

<< construct relation schemes from MVDs>>

do for each $m_i \in M$;

construct a relation scheme R_{i+m} consisting of all attributes appearing in m_i ; The set of all attributes that appears in m_i is the key of the relation scheme;

end;

$$R'' = \bigcup_{\substack{i=m+1\\i=m+1}}^{m+n} R_i;$$

Figure 3-2.

Figure 3-2 (continued) <<pre><<put all relation schemes together>> R = R' U R'';end PART-1; PART-2: <<this part removes a relation R if it has the same set of attributes as R_{i} , and the relationships of its attributes are not as strong as $R_{j}^{>>}$ <u>do</u> for each $R_i \in R$ for $1 \leq i \leq r - 1$; <<r is cardinality of R>> <u>do</u> for each $R_j \in R$ for $i + 1 \le j \le r$; $\underbrace{if}_{R_{i}} = U_{R_{i}} \underbrace{then}_{if} |K_{R_{i}}| \ge |K_{R_{i}}| \underbrace{then}_{if}$ $R = R - R_i; << K_{R_i} and K_{R_j}$ are <u>else</u> $R = R - R_j$; keys of R_i and R;>> end; end; end PART-2;

end algorithm;

(with key K_{R_i}) from the schema R, only if there exists a relation scheme R_j (with key K_{R_i}) in R such that:

(i)
$$R_i = R_i$$
,

or

(ii)
$$U_{R_i} = U_{R_j}$$
 and $|K_{R_i}| \ge |K_{R_j}|$.

The proof for the case of condition (i) is straightforward and needs no explanation. For the other case, suppose relation schemes R_i and R_j are constructed from dependencies $F_i: X \longrightarrow Y$ and $f_j: X' \longrightarrow Y'$ respectively. Since X U Y = X' U Y' (by assumption), and if we assume X' ^C X, then f_i is embodied in $R_j(X', Y')$ (which is the same as $R_i(X', \alpha, Y)$), and thus its information will not be lost if we delete R_i from the schema. A similar discussion now can be given if we relax on the assumption X' ^C X, completing the proof. #

Although the new algorithm eliminates some of the problems discussed earlier, the nonredundancy and separation principles cannot be satisfied. The following examples clarify these concepts.

Example 3.3: In our University Housing data base, if each COMPLEX can specify a set of RENTs for its various apartments, and if a COMPLEX together with the APT-CAPACITY can determine the RENT of the apartment, then we may have the following dependencies:

> $g_1: COMPLEX, APT-CAPACITY \longrightarrow RENT$ $m_1: COMPLEX \longrightarrow RENT$

and we can construct the following relation schemes:

R₁ (<u>COMPLEX</u>, <u>APT-CAPACITY</u>, RENT)

 R_2 (<u>COMPLEX</u>, <u>RENT</u>)

Although relation scheme R_2 of the above schema is redundant, but it is not removed by the Algorithm 3-2, and thus violating nonredundancy principle. Indeed, the schema does not satisfy the separation principle either, since it is not even in 2NF (it is easy to verify that in relation scheme R_1 , attribute RENT is partially dependent on the key). Example 3.4: Again in our University Housing data base, if each TENANT lives in only one apartment (APT#) which has one TYPE and is in one COMPLEX, and if each COMPLEX has only one TYPE of apartments, then we may have the following functional dependencies:

f,: TENANT -> APT#, COMPLEX, TYPE

 $f_2: COMPLEX \longrightarrow TYPE$

using the Algorithm 3-2, we can construct the following relation schemes:

R, (TENANT, APT#, COMPLEX, TYPE)

R₂ (<u>COMPLEX</u>, TYPE)

Notice that, in relation scheme R₁, attribute TYPE is transitively dependent on key TENANT, and thus violates the 3NF (although it is in 2NF).

The above examples indicate that the Algorithm 3-2 may produce a schema violating either nonredundancy principle or

separation principle (and in some cases both).

As a matter of fact, all of these problems arose mainly because we have completely ignored the composing rules for FDs and MVDs discussed in Chapter II. Thus, for further modifications of the algorithm these rules will be also considered.

3.5. The Dependency-to-Relation Method, 2nd Improvement

Considering the theoretical tools discussed in Chapter II, we can now further modify our algorithm. We will first modify the algorithm using the augmentation rule, and will prove that the schema constructed in this way is at least in second normal form (but not necessarily in 3NF or 4NF). Then we will have one more improvement to our method, this time considering the transitivity rule. We will see that these considerations eliminate many of the problems shown previously in our examples. Finally, we will present an algorithm considering all inference rules for both, FDs and MVDs. Then we will prove that the relational schema constructed by the algorithm satisfies all of the three design principles.

It is important to note, that in the algorithm which follows (Figure 3-3), only <u>explicit</u> dependencies (i.e., those FDs and MVDs given as part of the data base description) are considered, and therefore, not all of the problems may be eliminated by this algorithm. Thus, in Section 3.6 we will further modify the algorithm by considering not only explicit

FDs and MVDs, but also considering <u>implicit</u> dependencies (i.e., dependencies that are not explicitly presented as part of the data base description, but can be inferred from them by applying the inference rules) as well.

3.5.1. Proof of Representation

<u>Theorem 3.4</u>. Let R be a relational schema constructed from a set of dependencies $G = F \cup M$ using the Algorithm 3-3, then R represents the same information as G.

<u>Proof</u>: The proof for the schema constructed by the first two parts of the algorithm follows from Theorem 3.3. Third part of the algorithm may remove either an extraneous attribute from a relation scheme of the schema, or an entire relation scheme of the schema. We follow the proofs for both cases. (i). Removing extraneous attributes from a relation scheme: a set of attributes is removed from a relation scheme R_i only if it is a subset of non-key attributrs of another relation scheme R_j whose key is a subset of the key of relation R_i . Therefore, only those attributes are removed from a relation, but their relationship/s are also represented in a stronger form than original relation.

(ii). Removing a relation scheme from the schema: Relation scheme R_i is removed from the schema R only if after removing its extraneous attributes it becomes the same as another scheme in the schema. Thus its elimination from the schema

Algorithm 3-3

An Algorithm to Construct a Schema from a Set of

Dependencies; 2nd Improvement

INPUT: A set F of m FDs; and a set M of n MVDs.

OUTPUT: A set R of relation schemes (in lNF or better).

PART-1: <<this part constructs one relation scheme for each dependency>>

<< construct relation schemes from FDs>>

do for each f; E F;

construct a relation scheme R_i consisting of all attributes appearing in f_i ; The set of attributes that appears on the left side of f_i is the key of the relation scheme;

end;

$$\begin{aligned}
 \mathbf{R}' &= & \bigcup & \mathbf{R}_{i}; \\
 \mathbf{i} &= & \mathbf{i}
 \end{aligned}$$

<<construct relation schemes from MVDs>>

do for each $m_i \in M$;

construct a relation scheme R_{i+m} consisting of all attributes appearing in m_i ; The set of all attributes that appears in m_i is the key of the relation scheme;

end;

m+n R" = U i=m+1

Figure 3-3.

Figure 3-3 (continued)

<<pre><<put all relation schemes together>>

R = R' U R'';

end PART-1;

PART-2: <<this part removes a relation R_i if it has the same set of attribures as R_j, and the relationships of its attributes are not as strong as

<u>do</u> for each $R_i \in R$ for $1 \le i \le r - 1$; <<r is cardinality of R>>

 $\frac{\text{do for each } R_{j} \in R \text{ for } i + 1 \leq j \leq r;}{\text{if } U_{R_{j}} = U_{R_{j}} \underbrace{\text{then } if}_{R_{j}} | K_{R_{j}} | \geq |K_{R_{j}}| \underbrace{\text{then }}_{R_{j}} | \underbrace{\text{then }}_{R_{j}} | R = R - R_{i}; << K_{R_{j}} and K_{R_{j}} are \underbrace{\text{else }}_{R = R - R_{j}; keys of R_{i} and R_{j} >> keys of R_{i} and R_{i} > keys of R_{i} and R_{i} and R_{i} > keys of R_{i} and R_{i} a$

end;

end;

end PART-2;

PART-3: <<this part eliminates explicit partial dependencies
 from relation schemes>>

<u>do for</u> each $R_i \in R'$;

$$\frac{do \ for \ each \ R_{j} \ \varepsilon \ R' - R_{i}}{if \ K_{R_{j}} \ c \ K_{R_{i}} \ \frac{then}{then} [}$$

Figure 3-3 (continued)

$$\Psi = (U_{R_{i}} - K_{R_{i}}) \cap (U_{R_{j}} - K_{R_{j}}),$$

$$U_{R_{i}} = U_{R_{i}} - \Psi,$$

$$\underline{if} \quad U_{R_{i}} = K_{R_{i}} \quad \underline{then} \ R = (R - R_{i}) \cup R_{i}];$$

end;

end;

end PART-3;

.

•

end algorithm;

can be done without losing any information.

As the result, neither elimination of the extraneous attributes nor removal of the entire relation schemes (in PART-3), cause losing any information, completing the proof. #

Although the representation principle is not violated by the Algorithm 3-3 (Theorem 3.4), but the schema constructed in this way may not even be in 2NF, and thus violating separation principle. The following examples and discussions clarify this problem.

Example 3.5: Given the following dependencies:

 $f_1: A, B \longrightarrow C, D$ $f_2: B \longrightarrow D, E$

the algorithm first constructs the following relation schemes:

 $R_1(\underline{A}, \underline{B}, C, D)$ $R_2(\underline{B}, D, E)$

and then removes the extraneous attribute D from R_1 .

 $R_1(\underline{A}, \underline{B}, C)$ $R_2(\underline{B}, D, E)$

Although the schema of the above example is in 2NF (at least), but there are cases that the schema constructed by the Algorithm 3-3 is not necessarily in 2NF. This can simply be indicated by the following example. Example 3.6: Given the following dependencies:

 $f_1: A, B \longrightarrow C, D, E$

$$\begin{array}{ccc} f_2 \colon & B \longrightarrow & F \\ f_3 \colon & F \longrightarrow & D,G \\ f_4 \colon & A \longrightarrow & E \end{array}$$

the algorithm first constructs one relation scheme for each dependency as follows:

$$R_{1}(\underline{A},\underline{B},C,D,E)$$

$$R_{2}(\underline{B},F)$$

$$R_{3}(\underline{F},D,G)$$

$$R_{4}(\underline{A},E)$$

and then removes the extraneous attribute E from R_1 .

 $R_{1} (\underline{A}, \underline{B}, C, D)$ $R_{2} (\underline{B}, F)$ $R_{3} (\underline{F}, D, G)$ $R_{4} (\underline{A}, E)$

In the above schema, relation schemes R_2 , R_3 , and R_4 are in 2NF (at least). But how about R_1 ? R_1 is not in 2NF because, attribute D is partially dependent on key {A,B} (from f_2 and f_3 , and by Axiom FD3 we get FD B \rightarrow D,G; and from FD B \rightarrow D,G and Axiom FD6 we get FD B \rightarrow D).

The above problems arose mainly because we have considered only those dependencies that are explicitly given, and not those implied dependencies that can be inferred from the original set by applying inference rules. This is considered in our next step of modifications.

3.6. The Dependency-to-Relation Method, 3rd Improvement

Recall the Closure II and Closure III of a set of data dependencies discussed in Chapter II. We need to consider them here for further improvements to our algorithm. In addition, some subsets of these closures that are formally stated in the following definitions, are also of particular interest for further modifications.

<u>Definition 3.1</u>. Let $G = F \cup M$ be a set of mixed dependencies (FDs and MVDs), then <u>FD-Closure II</u> of G, denoted by FD-G^{II} is defined to be the set of all functional dependencies in G^{II} .

<u>Definition 3.2</u>. Let G F M be a set of mixed dependencies (FDs and MVDs), then <u>FD-Closure III</u> of G, denoted by $FD-G^{III}$ is defined to be the set of all functional dependencies in G^{III} .

Considering the FD-Closure III of data dependencies, now we can present the Algorithm 3-4 (Figure 3-4) which as we will prove later produces schema at least in 2NF.

3.6.1. Proof of Representation

<u>Theorem 3.5</u>. Let R be a relational schema constructed from a set of dependencies $G = F \cup M$ using the Algorithm 3-4, then R represents the same information as G.

<u>Proof</u>: From Theorem 3.4, it follows that the schema constructed by the first three parts of the algorithm embodies F and M. PART-4 of the algorithm removes a set of attributes Ψ from a relation R_i only if (i) Ψ is in FD-G^{III}, which means that it is represented (implicitly) by some other relation schemes, and (ii) the relationship of Ψ with other attributes in those relations is stronger than the relationship in R_i. In other

Algorithm 3-4

An Algorithm to Construct a Schema From a Set of

Dependencies; 3rd Improvement

A set F of m FDs; a set M of n MVDs; and FD-(F \cup M)^{III}. INPUT:

OUTPUT: A set R of relation schemes (in 2NF or better).

PART-1: <<this part constructs one relation scheme for each dependency>>

<< construct relation schemes from FDs>>

```
do for each f_i \in F;
```

construct a relation scheme ${\rm R}^{}_i$ consisting of all attributes appearing in f;; The set of attributes that appears on the left side of f_i is the key of the relation scheme;

.....

end;

 $\mathbf{R'} = \bigcup_{i=1}^{R} \mathbf{i};$ u la recenza de entre encontra estre estre ---and the second secon • - -<<construct relation schemes from MVDs>> do for each $m_i \in M$; المحاجمة المحاجم المحملة محملة المحملة محملة المحملة المحملة المحملة محملة construct a relation scheme R_{i+m} consisting of all attributes appearing in m_i; The set of all attributes that appears in m, is the key of the relation scheme; end; m+n R'' = U R;;i=m+1 anda Alaman ar
Figure 3-4 (continued)
<

```
<

```
relation schemes together>>
 R = R' U R";
end PART-1;
PART-2: <R_i if it has the
 same set of attributes as
$$R_j$$
, and the relationships
 of its attributes are not as strong as R_j >>
 do for each $R_i \in R$ for $1 \le i \le r - 1$; <>
 do for each $R_j \in R$ for $i + 1 \le j \le r$;
 if $U_{R_i} = U_{R_j}$ then if $|K_{R_i}| \ge |K_{R_j}|$ then
 $R = R - R_i$; <R_i and KR_j are
 else $R = R - R_j$; keys of R_i and
 end;
 Red;
 Red;
 Red = R - R_j; keys of R_i and
 Red
 Red = R - R_j; keys of R_i and
 Red
 Red = R_i + R_j + R
```
```

end PART-2;

PART-3: <<this part eliminates explicit partial dependencies
 from relation schemes>>

 $\frac{do \text{ for each } R_{i} \in R';}{do \text{ for each } R_{j} \in R; - R_{i};}$ $\frac{if K_{R_{j}} \subset K_{R_{i}} \quad \underline{then} [$ $\Psi = (U_{R_{i}} - K_{R_{i}}) \cap (U_{R_{j}} - K_{R_{j}}),$ $U_{R_{i}} = U_{R_{i}} - \Psi,$ $\frac{if}{I} \quad U_{R_{i}} = K_{R_{i}} \quad \underline{then} R = (R - R_{i}) \cup R_{i}];$

end;

end;

end PART-3;

PART-4: <<this part eliminates implicit partial dependencies from relation schemes>>

Trom retuction bon

repeat:

$$\underline{do \ for \ each \ g_{j} \ \varepsilon \ (FD-G^{III} - G);}$$

$$\underline{do \ for \ each \ R_{i} \ \varepsilon \ R';}$$

$$\underline{if \ LEFT}_{gj} \ C \ K_{R_{i}} \ \underline{then} \ [$$

$$\Psi = (U_{R_{i}} - K_{R_{i}}) \ \cap \ RIGHT_{gj},$$

$$U_{R_{i}} = U_{R_{i}} - \Psi,$$

$$\underline{if} \ U_{R_{i}} = K_{R_{i}} \ \underline{then} \ R = (R - R_{i}) \ U \ R_{i},$$

$$\underline{if} \ \Psi = \emptyset \ \underline{then} \ R' = R' \ U \ R_{h}(\omega), \ where \ \omega = (LEFT_{gj},$$

$$U \ \Psi)$$

and LEFT g is the key];

end;

end;

until no relation scheme is added to R';

end PART-4;

end algorithm;

words, only those attributes are removed from a relation scheme that not only are represented in other relation schemes of the schema, but their relationships are also represented (in a stronger form than in original relation) by them. Hence, in PART-4 of the algorithm no information would be lost, completing the proof. #

3.6.2. Proof of Normalization

The schema constructed by the Algorithm 3-4 is not in the desirable form (i.e., 4NF) yet, but it is in a better form than the schemas constructed previously. This claim is formalized by the following theorem.

Theorem 3.6. Let R be a relational schema constructed from a set of dependencies using the Algorithm 3-4, then R is in second normal form.

<u>Proof</u>: From Theorem 3.2, it follows that the schema constructed by the first part of the algorithm is in first normal form. Parts 3 and 4 of the algorithm eliminate all partial dependencies (both, explicit and implicit) from relation schemes, and thus leaving a 2NF schema. #

Algorithm 3-4 eliminates some of the problems discussed earlier, but not all of them. It produces a schema which is at least in 2NF, while the former algorithms could only guarantee the 1NF schemas. As it was stated in Chapter I, although a schema in 2NF has less undesirable properties than if it was in 1NF, but it is not yet completely free of all

62

undesirable properties, and thus, still violates the separation principle. The following examples clarify these concepts. example 3.7: Given the following dependencies:

$$f_1: A, D \longrightarrow B, C$$
$$f_2: A \longrightarrow B, E$$
$$m_1: A \longrightarrow D$$

the algorithm first constructs the following schemes:

$$R_{1}(\underline{A}, \underline{D}, B, C)$$
$$R_{2}(\underline{A}, B, E)$$
$$R_{3}(\underline{A}, \underline{D})$$

and then removes attribute B from R_1 .

$$R_{1}(\underline{A}, \underline{D}, C)$$

$$R_{2}(\underline{A}, B, E)$$

$$R_{3}(\underline{A}, \underline{D})$$

and since FD A \longrightarrow C holds in R₁ (see Figure 3-5), then attribute C is also partially dependent on key (in R₁), and thus is removed from R₁. Now, R₁(<u>A</u>, <u>D</u>) becomes the same as relation R₃ and hence it is removed. Finally, relation scheme R₄(<u>A</u>, C) is added to the schema by the algorithm.

 $R_{2}(\underline{A}, B, E)$ $R_{3}(\underline{A}, \underline{D})$ $R_{4}(\underline{A}, C)$

As the reader will notice, all partial dependencies are removed from the schemes, allowing the schema to be at Derivation of A \longrightarrow C from A,D \longrightarrow B,C and A \longrightarrow D.

Figure 3-5.

least in 2NF (in fact, the above schema is in 4NF).
Example 3.8: Given the following dependencies:

the algorithm produces the following two schemes:

$$R_1(\underline{A}, B, C, D)$$

 $R_2(\underline{B}, C)$

As we can see, the schema is in 2NF but not in 3NF. This is becuase, in relation scheme R_1 attribute C is transitivly dependent on key A (it is dependent on attribute B). This problem can be overcome by decomposing this type of relation schemes into their projections, using the transitive dependencies. This issue will be considered in the following section.

3.7. The Dependency-to-Relation Method, 4th Improvement

In this section we will further modify our algorithm to construct a schema at least in 3NF. To do so, we need to detect those dependencies that exist between two sets of non-prime attributes (or between a set of non-prime attributes and a proper subset of attributes making the key) in a relation scheme. This enables us to decompose the relation scheme into two of its projections, and thus eliminating many undesirable properties (if not all of them) from the schema. This is performed in PART-5 of the Algorithm 3-5 (Figure 3-6).

Algorithm 3-5

An Algorithm to Construct a Schema From a Set of

Dependencies; 4th Improvement

INPUT: A set F of m FDs; a set M of n MVDs; and FD-(F U M)^{III}.

OUTPUT: A set R of relation schemes (in 3NF or better).

<< construct relation schemes from FDs>>

<u>do</u> for each $f_i \in F$;

construct a relation scheme R_i consisting of all attributes appearing in f_i ; The set of attributes that appears on the left side of f_i is the key of the relation scheme;

end;

$$R' = \bigcup_{\substack{i=1}}^{m} R_{i};$$

<< construct relation schemes from MVDs>>

```
<u>do</u> for each m_i \in M;
```

construct a relation scheme R_{i+m} consisting of all attributes appearing in m_i ; The set of all attributes that appears in m_i is the key of the relation scheme:

end;

$$R'' = \bigcup_{\substack{n \in \mathbb{R} \\ i=m+1}}^{m+n} R;$$

<<pre><<put all relation schemes together>>

R = R' U R'';

end PART-1;

PART-2: <<this part removes a relation R_i if it has the same set of attributes as R_j, and the relation- ships of its attributes are not as strong as R_j>> <u>do for</u> each R_i ε R for 1 ≤ i ≤ r - 1; <<r is cardinality of R>>

$$\frac{do for each R_{j} \in R \text{ for } i + 1 \leq j \leq r;}{if U_{R_{j}} = U_{R_{j}} \frac{then}{if} |K_{R_{i}}| \geq |K_{R_{j}}| \frac{then}{R_{j}}}{R = R - R_{i}; << K_{R_{i}} \text{ and } K_{R_{j}}}{are}$$

$$\frac{else}{R} = R = R = R_{j}; \text{ keys of } R_{i} \text{ and}$$

$$\frac{end}{R_{j}} = \frac{R_{j}}{R_{j}} = \frac{R_{j$$

end;

end PART-2;

PART-3: <<this part eliminates explicit partial dependencies

from relation schemes>>

 $\underline{do \ for \ each \ R_{i} \ \varepsilon \ R';}$ $\underline{do \ for \ each \ R_{j} \ \varepsilon \ R' - R_{i};}$ $\underline{if \ K_{R_{j}} \ c \ K_{R_{i}} \ \underline{then} \ [}$ $\Psi = (U_{R_{i}} - K_{R_{i}}) \ \cap (U_{R_{j}} - K_{R_{j}}),$ $U_{R_{i}} = U_{R_{i}} - \Psi,$ $\underline{if \ U_{R_{i}} = K_{R_{i}} \ \underline{then} \ R = (R - R_{i}) \ U \ R_{i}];$

end;

end;

end PART-3;

PART-4: <<this part eliminates implicit partial dependencies from relation schemes>>

repeat:

$$\underline{do \ for \ each \ g_{j} \ \varepsilon \ (FD-G^{III} - G);}$$

$$\underline{do \ for \ each \ R_{i} \ \varepsilon \ R';}$$

$$\underline{if \ LEFT}_{g_{j}} \stackrel{c}{=} \ K_{R_{i}} \ \underline{then} \ [$$

$$\Psi = (U_{R_{i}} - K_{R_{i}}) \cap RIGHT_{g_{j}},$$

$$U_{R_{i}} = U_{R_{i}} - \Psi,$$

$$\underline{if} \ U_{R_{i}} = K_{R_{i}} \ \underline{then} \ R = (R - R_{i}) \cup R_{i},$$

$$\underline{if} \ \Psi = \emptyset \ \underline{then} \ R' = R' \cup R_{h}(\omega), \text{ where } \omega = LEFT_{g_{j}}$$

$$U \ \Psi)$$

and LEFT is gj the key];

end;

end;

until no relation scheme is added to R';

end PART-4;

PART-5: <<this part eliminates transitive dependencies

from relation schemes>>

repeat:

<u>do</u> for each $R_i \in R'$;

$$\begin{split} \Omega &= U_{R_{i}} - K_{R_{i}} <<\Omega \text{ is the set of non-key attributes}>>} \\ \underline{do \ for \ each \ g_{j} \ \varepsilon \ FD-G^{III};} \\ & \underline{if \ U_{g_{j}} \ c} \ \Omega \ \underline{then} \ [decompose \ R_{i} \ into \\ & R_{i_{1}} (U_{g_{j}}) \ \text{ with } LEFT_{g_{j}} \ as \ the \ key, \\ and \ & R_{i_{2}} (K_{R_{i}} \cup (\Omega-RIGHT_{g_{j}})) \ \text{ with } K_{R_{i}} \ as \ the \ key, \\ & R' = (R - R_{i}) \cup R_{i_{1}} \cup R_{i_{2}}; \\ end; \end{split}$$

•

end;

until no relation scheme is added to R':

end PART-5;

end algorithm;

.

3.7.1. Proof of Representation

Existence of an MVD in a relation is a necessary and sufficient condition for that relation to be decomposable into two of its projections without loss of information. That is, the original relation can be reconstructed by joining (natural join) those two projections. This is proved by Fagin in [18]. He has presented the following theorem.

<u>Theorem 3.7</u>. MVD X \longrightarrow Y holds for relation R(X,Y,Z) if and only if R is the join of its projections $R_1(X,Y)$ and $R_2(X,Z)$. <u>Proof</u>: It is simple to verify that R(X,Y,Z) is the join of its projections $R_1(X,Y)$ and $R_2(X,Z)$ if and only if the following condition holds: Whenever (x,y,z) and (x,y',z') are tuples of R, then so are (x,y',z) and (x,y,z'). This latter condition holds iff $Y_{XZ} = Y_{XZ}'$. The proof now follows from the definition of MVDs. #

From Fagin's theorem we can conclude the following corollary.

<u>Corollary 3.1</u>. The existence of a data dependency, FD or MVD, in a relation is a necessary and sufficient condition for that relation to be decomposable to its projections without loss of information.

<u>Proof</u>: The proof follows from Theorem 3.7 and the axiom FD-MVD1 which is: if an FD X \longrightarrow Y holds for R, then so does MVD X \longrightarrow Y. #

Considering Theorem 3.7 and Corollary 3.1, now we can

70

have the following theorem.

<u>Theorem 3.8.</u> Let R be a relational schema constructed from a set of dependencies $G = F \cup M$ using the Algorithm 3-5, then R represents the same information as G.

<u>Proof</u>: From Theorem 3.5 it follows that the relational schema constructed by the first four parts of the algorithm represents F and M. Fifth part of the algorithm decomposes a relation R_i into two of its projections R_i and R_i based on i_1 i_2 an FD (or an MVD) holding in it. Therefore, from Corollary 3.1 it follows that the original relation R_i can be reproduced from R_i and R_i using the natural join. Thus, R_i i_1 and R_i represent the same information represented by R_i completing the proof. #

3.7.2. Proof of Normalization

<u>Proposition 3.1</u>. If a relational schema is in nth Normal Form (nNF) for $2 \le n \le 4$, then the projection operation on relations never can convert the schema into a worse Normal Form, that is mNF for m < n.

Proof: The proof follows from the definition of Normal Forms
(separation principle), and from the definition of projection
operation. #

<u>Theorem 3.9</u>. Let R be a relational schema constructed from a set of dependenceis using the Algorithm 3-5, then R is in Third normal form.

Proof: From Theorem 3.6 it follows that the schema

71

constructed by the first four parts of the algorithm is in 2NF. And from Proposition 3.0 it follows that the fifth part of the algorithm cannot worsen the normalized form of the schema, thus leaving it to be still at least in 2NF. Now, because the schema is in 2NF, and since no transitive dependency exists in any relation scheme (those dependencies are all removed by PART-5), thus the schema is in 3NF, completing the proof. #

While Algorithm 3-5 eliminate all of the problems. The following examples clarify these concepts.

Example 3.9: Given the following FDs:

 $\begin{array}{rcl} f_1 \colon & A \longrightarrow & B, D, E, F \\ f_2 \colon & B \longrightarrow & C \\ f_3 \colon & C \longrightarrow & D \\ f_4 \colon & E \longrightarrow & F \end{array}$

the algorithm first constructs the following relation schemes: $R_1(\underline{A}, B, D, E, F)$ $R_2(\underline{B}, C)$ $R_3(\underline{C}, D)$ $R_4(\underline{E}, F)$ then it decomposes R_1 into $R_{11}(\underline{E}, F)$ and $R_{12}(\underline{A}, B, D, E)$, and

removes R_{11} since it is the same as R_4 .

$$\begin{array}{c} R_{12} (\underline{A}, B, D, E) \\ R_{2} (\underline{B}, C) \\ R_{3} (\underline{C}, D) \\ R_{4} (\underline{E}, F) \end{array}$$

and finally, since FD B \rightarrow D is in FD-G^{III} (it can be derived from f₂ and f₃ by applying axiom FD3), the algorithm decomposes R₁₂ into R₁₂₁(<u>B</u>,D) and R₁₂₂(<u>A</u>,B,E):

 $R_{121}(\underline{B}, D)$ $R_{122}(\underline{A}, B, E)$ $R_{2}(\underline{B}, C)$ $R_{3}(\underline{C}, D)$ $R_{4}(\underline{E}, F)$

As we can see the schema is in 3NF (it is in fact, in 4NF). Example 3.10: Given the following dependencies:

the algorithm first constructs the following relations:

 $R_1(\underline{A}, B, C, D)$

 $R_2(\underline{D},\underline{B})$

then, it decomposes R_1 into $R_{11}(\underline{D},B)$ and $R_{12}(\underline{A},C,D)$ because FD D \longrightarrow B holds in R_1 (from FD A \longrightarrow B,C,D and Axiom FD6 we can get FD A \longrightarrow B; from MVD D \longrightarrow B and FD A \longrightarrow B and Axiom FD-MVD2 we can get FD D \longrightarrow B).

$$R_{11}^{(\underline{D},B)}$$
$$R_{12}^{(\underline{A},C,D)}$$
$$R_{2}^{(\underline{D},\underline{B})}$$

The only problem with this schema is, that the relation $R_2(\underline{D},\underline{B})$ is now redundant and must be removed. In fact, we have already removed this kind of redundant relation schemes in PART-2 of the algorithm. Thus, to eliminate the above problem we can apply PART-2 of the algorithm for any new relation scheme constructed by the other parts of the algorithm. This consideration is shown in our Final algorithm. Example 3.11: Given the following dependencies:

$$f_{1}: A \longrightarrow B, C, D, E$$

$$f_{2}: A, B \longrightarrow C$$

$$f_{3}: D \longrightarrow E$$

$$m_{1}: F \longrightarrow G, E$$

the algorithm first constructs the following relations:

$$R_{1}(\underline{A}, B, C, D, E)$$

$$R_{2}(\underline{A}, \underline{B}, C)$$

$$R_{3}(\underline{D}, E)$$

$$R_{4}(\underline{F}, \underline{G}, \underline{E})$$

then it removes attribute C from R₂

$$R_{1}(\underline{A}, B, C, D, E)$$

$$R_{2}(\underline{A}, \underline{B})$$

$$R_{3}(\underline{D}, E)$$

$$R_{4}(\underline{F}, \underline{G}, \underline{E})$$

now R_1 is decomposed into $R_{11}(\underline{D}, E)$ and $R_{12}(\underline{A}, B, C, D)$ of which R_{11} is deleted, because it is the same as R_3 .

 $R_{12}(\underline{A}, \underline{B}, \underline{C}, \underline{D})$ $R_{2}(\underline{A}, \underline{B})$ $R_{3}(\underline{D}, \underline{E})$ $R_{4}(\underline{F}, \underline{G}, \underline{E})$

Note, the problem with this schema is not only that $R_2(\underline{A},\underline{B})$ is redundant, but also $R_4(\underline{F},\underline{G},\underline{E})$ is not in 4NF (because, from FD f₃ and MVD m₁ and axiom FD-MVD2 it follows that $F \longrightarrow E$).

To overcome the above problems we need to further modify our algorithm. This, which will be in fact, the final modification regarding the construction of a "desirable" schema is stated in the following section.

3.8. The Dependency-to-Relation Method, Final Improvement

Now, we are in the final state of the iterative refinement of the algorithm. The algorithm presented in this section (Figure 3-7), constructs a schema which will be proved to be in 4NF (hereafter, we will call this algorithm the FNF Algorithm). All proofs concerning the FNF Algorithm are given in Chapter V. In that chapter, we will formally examine the properties of the schema produced by the FNF Algorithm, and will propose directions for constructing optimal schemas.

The proofs of representation and normalization (and also nonredundancy) for the FNF Algorithm will be given in Chapter V. Here, we only give an example to show that the problems discussed earlier are eliminated by the final algorithm. The data base of the following example is taken from [18].

Example 3.12: Suppose we have a university data base with

75

Algorithm 3-6

<u>An Algorithm to Construct a Fourth Normal Form Schema</u> <u>from a Set of Dependencies: (FNF Algorithm)</u> INPUT: A set F of m FDs; a set M of n MVDs; the set $(F \cup M)^{II}$; the set $(F \cup M)^{III}$; and the set FD- $(F \cup M)^{III}$. OUTPUT: A set R of relation schemes in 4NF. PART-1: <<this part constructs one relation scheme for each dependency>> <<construct relation schemes from FDs>> <u>do for</u> each f_i \in F; construct a relation scheme R_i consisting of all attributes appearing in f_i; The set of attributes that appears on the left side of f_i is the key of the relation scheme;

end;

$$\begin{aligned}
 \mathbf{R}^{i} &= \bigcup_{\substack{i=1\\i=1}}^{m} \mathbf{R}_{i}^{i};
 \end{aligned}$$

<< construct relation schemes from MVDs>>

<u>do</u> for each $m_i \in M$;

construct a relation scheme R_{i+m} consisting of all attributes appearing in m_i ; The set of all attributes that appears in m_i is the key of the relation scheme;

end;

.

- ·

$$\Psi = (U_{R_{i}} - K_{R_{i}}) \cap (U_{R_{j}} - K_{R_{j}}),$$

$$U_{R_{i}} = U_{R_{i}} - \Psi$$

$$\underbrace{if}_{R_{i}} = K_{R_{i}} \underbrace{then}_{R} = (R - R_{i}) \cup R_{i};$$

end;

end;

end PART-3;

PART-4: <<this part eliminates implicit partial dependencies
 from relation schemes>>

repeat:

$$\frac{do for each g_{j} \varepsilon (FD-G^{III} - G);}{do for each R_{i} \varepsilon R';}$$

$$\frac{if LEFT_{g_{j}} \subset K_{R_{i}} then [}{\Psi = (U_{R_{i}} - K_{R_{i}}) \cap RIGHT_{g_{j}}, U_{R_{i}} = U_{R_{i}} - \Psi,$$

$$\frac{if U_{R_{i}}}{if U_{R_{i}}} = K_{R_{i}} then R = (R - R_{i}) \cup R_{i},$$

$$\frac{if \Psi = \emptyset then R' - R' \cup R_{h}(\omega), where \omega = (LEFT_{g_{j}} \cup \Psi)$$

and LEFT is the
 g;
key];

end;

end;

until no relation scheme is added to R';

end PART-4;

PART-5: <<this part eliminates transitive dependencies

from relation schemes>>

repeat:

```
\frac{do \ for \ each \ R_{i} \ \varepsilon \ R';}{\Omega = U_{R_{i}} - K_{R_{i}}} << \Omega \ is \ the \ set \ of \ non-key \ attributes>> 

<math display="block">\frac{do \ for \ each \ g_{j} \ \varepsilon \ FD-G^{III};}{\frac{if \ U_{g_{j}} \ \subseteq \ \Omega \ then \ [decompose \ R_{i} \ into \ R_{i_{1}} (U_{g_{j}}) \ with \ LEFT \ g_{j} \ as \ the \ key, 

and \ R_{i_{2}} (K_{R_{i}} \ U \ (\Omega-RIGHT_{g_{j}})) \ with \ K_{R_{i}} \ as \ the \ key, 

R' = (R' - R_{i}) \ U \ R_{i_{1}} \ U \ R_{i_{1}} \ U \ R_{i_{1}} \ U \ R_{i_{2}} ;
```

end;

end;

until no relation scheme is added to R':

end PART-5;

PART-6: <<this part decomposes decomposable schemes

constructed from FDs>>

repeat:

```
 \frac{\text{do for each } R_{i} \in R'; }{\Omega = U_{R_{i}} - K_{R_{i}}, <<\Omega \text{ is the set of non-key attributes}>} \\ \frac{\text{do for each } g_{j} \in G^{\text{III}}; \\ \frac{\text{if LEFT}_{g_{j}} \leq \Omega \text{ and } \text{RIGHT}_{g_{j}} \cap K_{R_{i}} \frac{\text{then } [\text{decompose } R_{i}] \\ \text{into } R_{i} (\text{RIGHT}_{g_{j}} \cup \text{LEFT}_{g_{j}}) \text{ with } \text{LEFT}_{g_{j}} \text{ as the} \\ 1 \qquad g_{j} \qquad g_{j} \qquad g_{j} \qquad g_{j}  key if:
```

g_jεFD-G^{III}, and (RIGHT U LEFT) otherwise; R. (LEFT U (U - RIGHT)) with all attributes ${}^{2}g_{j}$ ${}^{R_{i}}$ and appearing in R. as the key; ¹2 $R' = (R' - R_i) \cup R_i \cup R_i_2$ end; end: until no relation scheme is added to R'; end PART-6; PART-7: <<this part decomposes decomposable schemes constructed from MVDs>> repeat: <u>do</u> for each $g_i \in G^{II}$; <u>do</u> for each $R_i \in R'';$ if (LEFT U RIGHT) C R then [decompose R into R_{i_1} (LEFT U RIGHT) with LEFT as the key if: g_{j} $g_{j} \in FD-G^{II}$, and (RIGHT U LEFT) otherwise; g_j $R_{i_2} (U_{R_i} - RIGHT_{g_j})$ with all attributes appearing and in R_i as the key; $R^{"} = (R^{"} - R_{i}) \cup R_{i_{1}} \cup R_{i_{2}};$ end; end;

until no relation scheme is added to R"; repeat PART-2 for the new set R = R' U R";

end PART-7;

PART-8: <<this part removes those redundant schemes that are represented (implicitly) by other schemes of the schema>>

<u>do</u> for each $R_i \in R;$ <u>do</u> for each $g_j \in (G^{II} - G);$ <u>if</u> (LEFT $g_j \cup RIGHT_{g_j} = U_{R_i}$ and LEFT $g_j = K_{R_i}$ then $R = R - R_i;$

end;

end;

end PART-8;

end algorithm;

attributes CLASS, SECTION, STUDENT, MAJOR, EXAM, YEAR, IN-STRUCTOR, RANK, SALARY, TEXT, DAY, and ROOM. If a given CLASS consists of several SECTIONS, each of which has one INSTRUCTOR and various STUDENTS. And if each CLASS has a set of TEXTs, which are used by all SECTIONS of the CLASS. Also assume, that each SECTION of a CLASS meets on various DAYS, and on a given DAY, it has one meeting ROOM. We also know that, each STUDENT has only one MAJOR, and one YEAR. A STUDENT in a given CLASS and SECTION has several EXAM scores. And if each INSTRUCTOR has one RANK and one SALARY, then we can have the following FDs and MVDs:

- $f_1: CLASS, SECTION \longrightarrow INSTRUCTOR$
- $f_2: CLASS, SECTION, DAY \longrightarrow ROOM$
- $f_2:$ STUDENT \longrightarrow MAJOR, YEAR
- f_A : INSTRUCTOR \longrightarrow RANK, SALARY
- m₁: CLASS, SECTION ----> STUDENT, MAJOR, EXAM, YEAR
- m₂: CLASS, SECTION ---->> INSTRUCTOR, RANK, SALARY
- m_2 : CLASS, SECTION \longrightarrow DAY, ROOM
- m_{λ} : CLASS \longrightarrow TEXT
- m_{5} : CLASS, SECTION, STUDENT \longrightarrow EXAM

the FNF Algorithm first produces the following relation schemes (one scheme for each dependency):

R₁ (<u>CLASS</u>, <u>SECTION</u>, INSTRUCTOR) R₂ (<u>CLASS</u>, <u>SECTION</u>, <u>DAY</u>, ROOM) R₃ (<u>STUDENT</u>, MAJOR, YEAR) R₄ (INSTRUCTOR, RANK, SALARY)

R₅ (<u>CLASS</u>, <u>SECTION</u>, <u>STUDENT</u>, <u>MAJOR</u>, <u>EXAM</u>, <u>YEAR</u>)

R₆ (<u>CLASS</u>, <u>SECTION</u>, <u>INSTRUCTOR</u>, <u>RANK</u>, <u>SALARY</u>)

R7 (CLASS, SECTION, DAY, ROOM)

 $R_{o}(\underline{CLASS}, \underline{TEXT})$

R9 (CLASS, SECTION, STUDENT, EXAM)

then it decomposes R_5 into

R₅₁ (<u>STUDENT</u>, MAJOR, YEAR) and

R₅₂ (<u>CLASS</u>, <u>SECTION</u>, <u>STUDENT</u>, <u>EXAM</u>)

and it decomposes R_6 into

R₆₁ (<u>CLASS</u>, <u>SECTION</u>, INSTRUCTOR) and

R₆₂ (<u>CLASS</u>, <u>SECTION</u>, <u>RANK</u>, <u>SALARY</u>)

and since R_{51} is the same as R_3 , R_{52} is the same as R_9 , and R_{61} is the same as R_1 , they are all removed from the schema. Relation scheme R_{62} is also removed from the schema because it is represented by R_1 and R_2 (in an stronger form). For the same reason, R_7 which is represented by R_2 in a stronger form, is deleted. Therefore, the final schema constructed by the algorithm is:

> R₁ (<u>CLASS</u>, <u>SECTION</u>, INSTRUCTOR) R₂ (<u>CLASS</u>, <u>SECTION</u>, <u>DAY</u>, ROOM) R₃ (<u>STUDENT</u>, MAJOR, YEAR) R₄ (<u>INSTRUCTOR</u>, RANK, SALARY) R₈ (<u>CLASS</u>, TEXT) R₉ (<u>CLASS</u>, <u>SECTION</u>, <u>STUDENT</u>, <u>EXAM</u>)

As reader will notice the constructed schema is in fourth normal form (the proof and related formal discussions are given in Chapter V).

3.9. Chapter Summary and Remarks

A new design approach for constructing relational schemas from dependencies (FDs and MVDs) designated by the data base administrator, has been proposed and investigated. After an iterative refinement of the approach, a final algorithm (FNF) which produces 4NF relation schemes has been presented. It has been proved that the schema constructed by this method satisfies the three design principles (the proofs concerning the FNF Algorithm will be given in Chapter V), and thus it is free from undesirable properties discussed in Chapter I.

CHAPTER IV

CONSTRUCTING CLOSURE II AND CLOSURE III

4.1 Introduction

In this chapter two algorithms are presented for constructing Closure II and Closure III of a given set of data dependencies. Also, detailed analyses of these algorithms are given, and investigations are made to indicate that the worst case (i.e., the case in which all data dependencies designated by the data base administrator are completely chained together) of each algorithm is unlikely to happen for real world applications. Anyhow, to maintain the consistency of the design approach, this special case must also be considered. An algorithm which detects this case and warns the data base administrator of the "bad" situation is given at the end of this chapter.

4.2 Description of the Closure II Algorithm

An FD f: $X \longrightarrow Y$ is in <u>canonical form</u> if and only if |Y| = 1, that is, the right side of the function consists of <u>one</u> attribute [8]. A set of FDs is in canonical form if and only if all of its member FDs are in canonical form. <u>Proposition 4.1</u>. For any set of FDs, there exists an equivalent set (i.e., a set representing the same information) of FDs in canonical form.

<u>Proof</u>: Any FD f: X \longrightarrow Y (Y = {Y₁, Y₂,...,Y_n}, where n > 1) can easily be converted to a set of FDs F = {f₁, f₂,..., f_n} in canonical form using Axiom FD6.

$$f_{1}: X \to Y_{1},$$

$$f_{2}: X \to Y_{2},$$

$$\vdots$$

$$f_{n}: X \to Y_{n};$$

completing the proof. #

Considering proposition 4.1, we will assume without loss of generality that the set of FDs used as input to our Closure II Algorithm (Figure 4-1) is in canonical form. But, a similar assumption cannot be made for the set of MVDs, simply because having an MVD m:X \rightarrow {Y₁,Y₂, ...,Y_n}, does not necessarily result to

$$\begin{array}{l} \mathbf{m}_1 \colon \mathbf{X} \longrightarrow \mathbf{Y}_1, \\ \mathbf{m}_2 \colon \mathbf{X} \longrightarrow \mathbf{Y}_2, \\ \vdots \\ \mathbf{m}_n \colon \mathbf{X} \longrightarrow \mathbf{Y}_n. \end{array}$$

We will also assume that for every FD X \rightarrow Y (or MVD X \rightarrow Y) used as input, Y $\not\subset$ X. This is a reasonable assumption for real world applications because, FD X \rightarrow X (or MVD X \rightarrow X) means, having a set of attributes X, we can determine (or multidetermine) a set of attributes X, and clearly, having a set X at our disposal we do not

Algorithm 4-1

An Algorithm to construct the Closure II of a set of data dependencies A set F of FDs in canonical form; a set M of MVDs; INPUT: and G = F || M. OUTPUT: Closure II of G, G^{II} ; FD- G^{II} ; and MV- G^{II} . $D = J^{\dagger} = F;$ V = J'' = M;repeat $H' = \emptyset;$ $H'' = \emptyset;$ do for each f; EF+M; <u>do</u> <u>for</u> each f_j_EJ'+J"; $\underline{if}_{RIGHT}_{f_i} \overset{\tilde{\subseteq}}{=} {}^{LEFT}_{f_i} \overset{\underline{then}}{=}$ [if $f_j \in J'$ [construct an FD f_k with $\operatorname{LEFT}_{f_k} = \operatorname{LEFT}_{i} (\operatorname{LEFT}_{f_i} - \operatorname{RIGHT}_{f_i})$ and RIGHT_f=RIGHT_f; $H' = H ||f_{k};$ else;]; else [construct an MVD f_k with $\text{LEFT}_{f_k} = \text{LEFT}_{f_i} \bigcup (\text{LEFT}_{f_i} = \text{RIGHT}_{f_i})$ and RIGHT f_k=RIGHT f_j-LEFT f_j; $\underline{if} \operatorname{RIGHT}_{f_k} \subseteq \operatorname{LEFT}_{f_k} \underline{then}$ $H'' = H''[f_k]];$

Figure 4-1

Figure 4-1 continued

,

.

.

have to determine (or multidetermine) it. Anyhow, those kind of data dependencies must be detected (if any) to avoid inconsistency of the approach. A detection algorithm is presented at the end of this chapter.

4.2.1 Proof of termination

Theorem 4.1. Algorithm 4-1 halts.

<u>Proof</u>: Both, inner loop and innermost loop of the algorithm are finite loops because F + M and J' + j'' are finite sets. At each iteration of outer loop, the innermost loop generates a number of new data dependencies. Since the cardinality of the set generated at pass i (for $2 \le i \le n$) is less than the cardinality of the set generated at pass i - 1 (because, for the worst case in which all dependencies are chained together, if p dependencies are generated at pass i - 1, the number of dependencies generated at pass i is p - 1), then the algorithm ultimately reaches to the point that $H' + H'' = \emptyset$, and hence it halts. #

4.2.2 Proof of correctness

<u>Theorem 4.2</u>. Upon termination of the Algorithm 4-1, G^{II} is the Closure II of G, FD- G^{II} contains all FDs in G^{II} , and MV- G^{II} contains all MVDs in G^{II} . <u>Proof</u>: As we can observe, G^{II} is calculated by the algorithm as $G^{II} = G + \bigcup_{i=1}^{n} (H' + H'')_{i}$, where $(H' + H'')_{i}$ is the set of dependencies generated at pass i. Having G at the first pass, the set $(H' + H'')_1$ generated, is in fact G^2 , and similarly set $(H' + H'')_{n-1}$ generated at pass n - 1 is G^n^{\dagger} . Therefore, we have $G^{II} = G \bigcup G^2 \bigcup G^3 \bigcup \ldots \bigcup G^n = \bigcup_{i=1}^n G^i$. Now, the proof of the correctness for this formula will be similar to the case that G^{II} is the transitive closure of G [28]. The proof is in two parts. 1) $\bigcup_{i=1}^n G^i \subset G^{II}$, where n is the smallest integer satisfying $G^n \bigcup G^{n+1} = G^n$. We prove this part by induction.

(basis)

From Definition 2.1 (given in Chapter II), it follows that $G \subset G^{II}$.

(induction step)

Suppose $G^{i} \subset G^{II}$, $i \ge 1$, and let $\langle A, B \rangle \in G^{i+1}$, where $\langle A, B \rangle$ means $A \longrightarrow B$ (or $A \longrightarrow B$). Since $G^{i+1} = G^{i}G$, there exists some $C \in U$ such that $\langle (A - C), (C - A) \rangle \in G^{i}$ and $\langle C, B \rangle$ ϵ G. By the induction hypothesis and the basis step, $\langle (A - C), (C - A) \rangle \in G^{II}$ and $C, B \in G^{II}$. Since G^{II} is closed under axioms FD4 and MVD4, it follows that $\langle A, B \rangle \in G^{II}$, thus completing the induction.

2) $G^{II} \subset \bigcup_{i=1}^{n} G^{i}$. Let <A,B> and <B,C> be arbitrary elements of $\bigcup_{i=1}^{n} G^{i}$. Then for some positive integers p and q, <A,B> εG^{p} and <B α ,C> εG^{q} . Then <A α ,C> $\varepsilon G^{p}G^{q}$ (or <A α ,C> εG^{p+q}).

t The concept is borrowed from the following definition: Def. Let R be a binary relation on a set A and let $n \in \mathbb{N}$. Then, \mathbb{R}^n is defined as follows:

1. R^0 is the relation of equality on the set A. 2. $R^{n+1} = R^n R$ Therefore, $\langle A\alpha, C \rangle_{i=1}^{n} G^{i}$ and hence $\prod_{i=1}^{n} G^{i}$ contains all dependencies derivable from G using Axioms FD4 and MVD4, completing the proof of the first part of the theorem.

We claim also that the type of data dependency (FD or MVD) is considered by the algorithm appropriately. In fact, for each pair of data dependencies f_i and f_j to be checked, four possibilities should be considered:

- Both f_i and f_j belong to the family of FDs. For this case the generated dependency would be of type FD.
- 2) f_i belongs to FDs family, and f_j belongs to MVDs family. Since f_i belongs to FDs family, then it also belongs to MVDs family (Axiom FD-MVD1). Thus, both f_i and f_j are MVDs, and the generated dependency is also an MVD.
- 3) Both f_i and f_j belong to MVDs family. As for the case₂ the generated data dependency would be of type MVD.
- 4) f_i belongs to MVDs family, and f_j belongs to FDs family. Again, since f_j is an FD, it is also an MVD, and thus we can generate a data dependency f_m of type MVD whose right side is the same with the right side of f_j . Now, by applying FD-MVD2 to the pair f_m and f_j we can generate an FD f_k , because (i) LEFT $f_j \subseteq LEFT f_m$ (in fact, we showed that LEFT $f_j = LEFT f_m$), and (ii) RIGHT $f_m \cap LEFT_{f_j} = \emptyset$ (by assumption). Therefore, the

resulting dependency for this case is of type FD.

From the above discussion we conclude that the generated data dependency is of type FD if and only if f_j belongs to the set of FDs, and it is an MVD otherwise. This is considered in innermost loop of the algorithm (the statement, if $f_j \in J' \dots$), thus completing the proof. #

4.2.3 Output analysis

At each iteration of the outer loop, a new data dependency (either FD or MVD) is generated if and only if there <u>is</u> a connection (in the sense of Axioms FD4 and MVD4) between a pair of dependencies. More connections between pairs of dependencies exists, more new data dependencies are generated. Thus, the best case is in fact when there is no connection between any pair, and therefore, output is only as large as input. Suppose there are m FDs and n MVDs as the input to the algorithm, then $|G^{II}|=|F|+|M|=$ m+n. On the other hand, the worst case happens when all given data dependencies are completely chained together. For this case, the number of data dependencies generated by the algorithm is proportional to $(m+n)^2$, which is proved by the following theorem.

<u>Theorem 4.3</u>. Let G^{II} be the Closure II of G = (F | M) constructed by the Algorithm 4-1, then for the worst case (all dependencies chained together), $|G^{II}| = (m+n)(m+n+1)/2$, where m = |F| and n = |M|.

92

Proof: By induction:

(m + n = 1)

If only one data dependency is involved, no new data dependency is generated by the algorithm (because the process requires at least a <u>pair</u> of dependencies), and thus the output consists of only 1 data dependency, that is, the original one. Using the above formula for m + n = 1, we also get $|G^{II}| = 1$.

(induction step)

If for i data dependencies $|G_i^{II}| = i \cdot (i+1)/2$, and if we add a new data dependency to our input set, and since by our assumption of the worst case, all i data dependencies have connections with the (i + 1)th data dependency, thus i more dependencies are generated. This has to be added to 1 (for the added dependency itself) to get the total number of data dependencies, $|G_{i+1}^{II}| =$ $i \cdot (i + 1)/2 + (i + 1) = (i + 1) \cdot ((i + 1) + 1)/2$, thus completing the induction. #

As we saw, the formula of Theorem 4.3, is for a situation in which <u>all</u> data dependencies designated by data base administrator are chained together. As an example, consider the following case where G is

 $f_{1}: A_{1} \rightarrow A_{2},$ $f_{2}: A_{2} \rightarrow A_{3},$ $f_{3}: A_{3} \rightarrow A_{4},$:

$$f_{i}: A_{i} \rightarrow A_{i+1}'$$

$$f_{m}: A_{m} \rightarrow A_{m+1},$$

$$m_{1}: A_{m+1} \rightarrow A_{m+2},$$

$$m_{2}: A_{m+2} \rightarrow A_{m+3},$$

$$\vdots$$

$$m_{j}: A_{m+j} \rightarrow A_{m+j+1},$$

$$\vdots$$

$$m_{n}: A_{m+n} \rightarrow A_{m+n+1}.$$

Now, suppose there is a "gap" (to cut the chain) somewhere in the chain, say, close to the middle of the set, then we have two subsets each with the approximate length of (m+n)/2 and both in the worst case (all elements of each subset are fully chained together). Then to find the number of elements in G^{II} , we can use the formula of Theorem 4.3 for each subset and then multiply the result by 2. Similarly, for 2 gaps in the set, we get 3 subsets, and we need to compute the number of data dependencies generated in one of them and multiply it by 3 (assuming that all subsets are of the same length). By the same token, if there are j gaps in the set, we get

$$|G^{II}| = \frac{(m+n)(m+n+j+1)}{2(j+1)}$$
 (4.2.2.1)

To have a better understanding of the formula (4.2.2.1) for the real world applications, we should notice that, if (m+n) is not very large (say < 100),

then even if J is small, $|G^{II}|$ won't be too large. And for the large value of (m+n), J will be comparatively large too (from the fact that not most of the dependencies in a <u>large</u> data base are fully chained), and again $|G^{II}|$ won't be too large.

4.2.4 Speed analysis

At the first pass through the outer loop, and for each iteration of the inner loop, the innermost loop executes (m+n) times (m=|F|, and n = |M|), and since inner loop itself is bounded by m + n iterations, then the entire loop executes (m+n)² times. If no new dependency is generated by this pass (i.e., the best case -no connections between any pair of dependencies), then algorithm terminates in time $0(p^2)$, where p = m+n. Considering the worst case, the outer loop iterates p times (at each iteration the number of dependencies generated is one dependency less than the previous iteration), and hence, the overall time for the algorithm is below $0(p^3)$.

> (best case) $T_{b} = p^{2}$ (worst case) $T_{w} = p \cdot (p) + p \cdot (p-1) + p \cdot (p-2) + \dots + 1 = p^{2} \cdot (p+1) / 2$

4.3 Description of the Closure III Algorithm

All assumptions for this algorithm are the same as what we had for the Algorithm 4-1. The basic objective of the algorithm is to construct G^{III} -- the

95
Closure III of a given set of data dependencies G, using Axioms FD4, MVD4, and FD-MVD2. Since G^{II} , which is the result of applying Axioms FD4 and MVD4 to the set G, can be constructed by the Algorithm 4-1, then the algorithm 4-2 gets as input G^{II} of G, and applies Axiom FD-MVD2 to G^{II} to get G^{III} (See Figure 4.2).

Example 4.1: This carefully selected example indicates the significance of the outer loop of the algorithm. The proofs of correctness and termination are given later.

Suppose the following set of data dependencies are given as the input to Algorithm 4-2:

$$\begin{split} \mathbf{m}_1 &: \mathbf{A} \longrightarrow \{\mathbf{B}, \mathbf{C}, \mathbf{D}\}, \\ \mathbf{m}_2 &: \mathbf{E} \longrightarrow \{\mathbf{B}, \mathbf{F}\}, \\ \mathbf{m}_3 &: \mathbf{G}, \mathbf{H} \longrightarrow \{\mathbf{B}, \mathbf{I}\}, \\ \mathbf{f}_1 &: \mathbf{F}, \mathbf{C}, \mathbf{E} \longrightarrow \{\mathbf{B}\}. \end{split}$$

then, at the first pass the algorithm generates the FD $f_{11}: G, H \rightarrow \{B\}$

and then using FD f11, another FD is generated

 $f_{12}: E \rightarrow \{B\}$

and finally, using FD f_{12} , the algorithm generates

$$f_{13}: A \rightarrow \{B\}.$$

4.3.1. Proof of correctness

<u>Theorem 4.4</u>. Upon termination of the Algorithm 4-2, all FDs derivable from the set G^{II} by applying FD-MVD2 are in FD- G^{III} , and thus G^{III} is the Closure III of G.

Algorithm 4-2 An Algorithm to construct the Closure III of a set of data dependencies INPUT : MV-G^{II} and FD-G^{II} of a set of data dependencies G; OUTPUT: G^{III} and FD-G^{III} of G; $V = MV - G^{II};$ $D = FD - G^{II};$ FD-G^{III} = FD-G^{II}; repeat $H = V^{t} = \emptyset;$ <u>do</u> for each $f_i \in V$; <u>do</u> for each $f_i \in D$; $\underline{if} \operatorname{RIGHT}_{f_i} \subseteq \operatorname{RIGHT}_{f_i} \underline{then}$ $[if RIGHT_{f_i}] \cap LEFT_{f_j} = \emptyset$ then [construct an FD f_k with $\text{LEFT}_{f_k} = \text{LEFT}_{f_i}$ and RIGHT_{FK} = RIGHT_{f;}; $\underbrace{\texttt{if RIGHT}}_{f_k} \underbrace{\not \subset}_{k} \underbrace{\texttt{LEFT}}_{f_k} \underbrace{\texttt{then}}_{k}$ $H = H \bigcup_{k}^{n} f_{k};$ else;] <u>else</u> $V' = V' [[f_i;]]$

end;

end;

. •

Figure 4-2

Figure 4-2 continued

.

$$V = V';$$

$$D = H;$$

$$FD-G^{III} = FD-G^{III} \bigcup H;$$

$$\underbrace{until}_{G} H = \emptyset \text{ or } V' = \emptyset;$$

$$G^{III} = FD-G^{III} + MV-G^{II};$$

end algorithm;

Proof: At the first pass of the algorithm, all MVDs in the set are checked with all FDs (pair wise) and the new FDs are generated. For the following passes, not all of the MVDs of the original set are needed to be checked with new FDs, and therefore, algorithm considers only those MVDs that have a chance of generating new FDs. This can be formalized as follows:

Suppose the input to the algorithm is

$$f_{1}: X_{1} \rightarrow Y_{1}$$

$$f_{2}: X_{2} \rightarrow Y_{2}$$

$$i$$

$$f_{j}: X_{j} \rightarrow Y_{j}$$

$$f_{m}: X_{m} \rightarrow Y_{m}$$

$$m_{1}: W_{1} \rightarrow Z_{1}$$

$$m_{2}: W_{2} \rightarrow Z_{2}$$

$$i$$

$$m_{n}: W_{n} \rightarrow Z_{n}$$

where sets X's, Y's, W's and Z's are not necessarily disjoint. At pass i, an MVD $W_j \longrightarrow Z_j$ is removed and will not be considered for pass i+1, if either

1)
$$Y_{i} \not \in Z_{j}$$
 for $1 \le i \le m$,
or

2) for all
$$y_i \in Z_j$$
 X_i and Z_j are disjoint --
 $X_i \cap Z_j = \emptyset$

Part (1): If the set H of new FDs generated by pass i is

$$\begin{array}{ll} f_{i1} \colon & W'_1 \to Y'_1 \\ f_{i2} \colon & W'_2 \to Y'_2 \\ \vdots \\ f_{ip} \colon & W'_p \to Y'_p, \end{array}$$

then the existence of MVD $W_j \rightarrow Z_j$ in the set of MVDs in pass i+1 is not significant, and it does not generate any newer Functional Dependency. In fact, if any FD is generated from MVD $W_j \rightarrow Z_j$ and the set H, then it must be an FD $W_j \rightarrow Y'_k$ (for some value of k in $1 \le k \le p$). This can happen only if $Y'_k \subseteq Z_j$ (Axiom FD-MVD2), and it is a contradiction to our assumption because we have

$$Y_k \in \{Y_1, Y_2, \dots, Y_p\}$$
 and

$$\{Y_1, Y_2, \dots, Y_p\} \subseteq \{Y_1, Y_2, \dots, Y_m\}$$

and thus

$$\mathbf{Y}_{\mathbf{k}} \in \{\mathbf{Y}_{1}, \mathbf{Y}_{2}, \dots, \mathbf{Y}_{m}\}$$

and by our assumption, no member of $\{Y_1, Y_2, \dots, Y_m\}$ can be a subset of Z_j .

Part (2): Again, to prove the redundancy of the MVD $W_j \rightarrow Z_j$ for pass i+1, we need to prove that its existence does not cause generating any newer FD. In fact, if MVD $W_j \rightarrow Z_j$ and the new set H of FDs, can generate a newer functional dependency, then it must be FD $W_j \rightarrow Y'_k$ (for some value of k in $1 \le k \le p$). This is because there must be an FD $W'_k \rightarrow Y'_k$ such that $Y'_k \subseteq Z_j$ and $W'_k \cap Z_j = \emptyset$. If $Y'_k \subseteq Z_j$, then since we have $\begin{array}{l} \mathbf{Y}_{k}^{*} \in \{\mathbf{Y}_{1}^{*}, \mathbf{Y}_{2}^{*}, \dots, \mathbf{Y}_{p}^{*}\} \\ \{\mathbf{Y}_{1}^{*}, \mathbf{Y}_{2}^{*}, \dots, \mathbf{Y}_{p}^{*}\} \subseteq \{\mathbf{Y}_{1}^{*}, \mathbf{Y}_{2}^{*}, \dots, \mathbf{Y}_{m}^{*}\} \end{array}$

then we get

 $Y'_k = Y_h$ for some value of h in $1 \le h \le m$.

And by our assumption, for all $Y_i \subseteq Z_j$ (including Y_h), we have $X_i \cap Z_j = \emptyset$. Thus MVD $W_j \longrightarrow Y'_k$ could be generated in pass i, and there is no need to carry MVD $W_j \longrightarrow Z_j$ to pass i+1, completing the proof. #

4.3.2 Proof of termination

Theorem 4.5 Algorithm 4-2 halts.

<u>Proof</u>: Both inner loop and innermost loop of the algorithm are finite loops because V and D are finite sets for all iterations of the outer loop. Since at each iteration of the outer loop either set H or set V (sometimes both) becomes smaller, then the algorithm always halts. #

4.3.3 Output Analysis

<u>Proposition 4.2</u> Let F be a set of FDs and M be a set of MVDs, then the maximum possible number of FDs that can be generated from F and M using Axion FD-MVD2 is $|F| \cdot |M|$. Proof: Consider the following data dependencies:

$$\begin{array}{ccc} \mathbf{f}_1 \colon & \mathbf{X}_1 \to \mathbf{Y}_1 \\ \mathbf{f}_2 \colon & \mathbf{X}_2 \to \mathbf{Y}_2 \\ \vdots \end{array}$$

$$f_{m}: X_{m} \rightarrow Y_{m}$$

$$m_{1}: W_{1} \rightarrow Z_{1}$$

$$m_{2}: W_{2} \rightarrow Z_{2}$$

$$\vdots$$

$$m_{n}: W_{n} \rightarrow Z_{n}$$

Any FD W' \rightarrow Y', generated from the above set by applying Axiom FD-MVD2 must have the following properties:

W'
$$\varepsilon \{W_1, W_2, \dots, W_n\}$$
 and
Y' $\varepsilon \{Y_1, Y_2, \dots, Y_m\}$.

Now, since there are n possibilities for W' and m possibilities for Y', then there are m.n (i.e., $|F| \cdot |M|$) possibilities for the FD W' \rightarrow Y'. #

4.3.4 Speed analysis

At each iteration of the inner loop, innermost loop executes (at most) m times (m = $|\text{FD-G}^{II}|$). At each iteration of the outer loop, inner loop iterates n times (n = $|\text{MV-G}^{II}|$). Therefore, at each iteration of the outer loop, the innermost loop executes m.n times. Considering proposition 4.2, and assuming that each single new functional dependency is generated at one iteration of the outer loop (the worst case), then the outer loop iterates m.n times and thus the overall time for the algorithm is $0(m.n)^2$. For the best case which most likely happens for real world applications, all expected FDs are generated at the first pass of the algorithm, and thus the algorithm halts in time O(m.n).

4.4 Description of the Detection Algorithms

As we mentioned earlier, there are some special cases that will never happen for the real world applications, when the design issue of the relational data base is of interest. Although we believe in this fact, we still care about the consistency of the approach, and thus those special cases will be detected (if they happen). Two algorithms that detect the two most unusual cases (in the sense of data dependencies designated by the data base administrator), and warn the data base administrator of the "bad" situations are presented in this section. Algorithm 4-3 (Figure 4-3) discovers the "fully chained data dependencies" situation, and Algorithm 4-4 (Figure 4-4) detects the "reflexive" data dependencies.

4.5 Chapter Summary and Remarks

Two algorithms for computing Closure II and Closure III of a given set of data dependencies have been presented. It has been shown that the worst case of these algorithms is unlikely to happen for real world applications. In spite of this fact, to maintain the consistency of the method, two algorithms have been presented to detect these worst cases (if they happen) and warn the data base administrator of the unexpected situation.

Algorithm 4-3

An Algorithm to detect the fully chained data dependencies situation

- INPUT : A set F of m FDs; and a set M of n MVDs;
- OUTPUT: "YES", if all data dependencies in (F \bigcup M) are chained together, and "NO" otherwise;
- STEP 1: (Initialization) m = k = i = 1; count = 0;
- STEP 2: (Get the right side of dependency) right=RIGHT_{gi}; j = 1;
- STEP 3; (Check if all dependencies are chained) if count = n - 1 then print "YES", STOP;
- STEP 4: (Check if this dependency has been tested) if i =
 j or MARK(i) = 1 then go to STEP 7;
- STEP 5: (Check the left side of dependency) if right q LEFT_{gj} or MARK(j) = 1 then go to STEP 7; PUSH [right,j,i] into STACK; MARK(i) = 1; count = count + 1; i = j; go to STEP 2;
- STEP 6: (Pop from stack). <u>if</u> STACK empty <u>go</u> to STEP 8; POP STACK in [right,j,i]; MARK(i)=0; count = count - 1;
- STEP 7: (Check if all dependencies have been tested for this right side) <u>if</u> j < n <u>then</u> j = j + 1, go to STEP 3; <u>else</u> go to STEP 6;
- STEP 8: (Termination) if m < n then m = m + 1, i = m, go
 to STEP 2; else print "NO", STOP;</pre>

Figure 4-3

Algorithm 4-4

An Algorithm to detect the reflexive data dependencies

INPUT : A set F of m FDs; and a set M of n MVDs; OUTPUT: A set $X \subseteq (F | J M)$ of reflexive data dependencies; $X = \emptyset;$ <u>do for</u> each $f_i \in F + M;$

$$\frac{\text{if } \text{RIGHT}_{f_i} \subseteq \text{LEFT}_{f_i}}{X = X \bigcup f_i};$$

end;

end algorithm;

Figure 4-4

CHAPTER V

ANALYSIS OF RELATIONS BY DEPENDENCY-TO-RELATION

5.1 Introduction

In Chapter III an algorithm (FNF Algorithm) was presented for constructing a relational schema from a set of functional dependencies and multivalued dependencies. In this chapter properties of the relation schemes constructed by the FNF Algorithm are discussed (section 5.2), and it is shown that although the schema constructed is desirable (i.e., contains no undesirable properties), it is not necessarily "optimal". Heuristics for choosing "reasonable" decompositions which lead to an "optimal" schema are suggested, and an algorithm which decomposes a set of decomposable schemes into an optimal form is given. Finally, two previous approaches are examined in sections 5.3 and 5.4, and they are compared with the new approach.

5.2 Properties of Relations constructed by FNF Algorithm

Recall the three design principles (Representation, Separation, and Nonredundancy) given in Chapter II. It is important to characterize the properties of the relational schema constructed by the FNF Algorithm to see if all three design principles are satisfied, and thus, to determine whether the schema is an "ideal" schema or not.

5.2.1 Representation Principle

The schema constructed by our FNF Algorithm, clearly, must not violate the representation principle. It must represent all information of interest (i.e., the same information as input). Different researchers have defined the concept of "the same information" in different ways. Four different definitions are discussed and compared in [6]. In this section, we first describe the most common and reasonable definition, and we will show that the schema constructed by the FNF Algorithm satisfies the representation principle (assuming this definition). Then representation is defined in another way, a definition which our schema of FNF Algorithm does not always satisfy. These alternatives are studied, and suggestions are given to the possible resolution of the paradox.

<u>Definition 5.1</u>. Let R and S be two relational schemas, then S represents the same information as R if they have the same attributes and databases of S contain the same data as the databases of R.

The above definition which relys on the corollary 3.1 given in Chapter III, can be clarified by the following examples.

Example 5.1: Suppose an schema R consists of a relation scheme TACT (TENANT, ATP#, COMPLEX, TYPE), and in addition to FD TENANT \rightarrow APT#, COMPLEX, TYPE, the FD COMPLEX \rightarrow TYPE holds in TACT, and suppose an instance of relation TACT is as in Figure 5-1a. Now, if the schema S consists of two relation

schemes CT (COMPLEX, TYPE) and TAC (TENANT, APT#, COMPLEX) which are indeed, the projections of R on COMPLEX and TYPE; and on TENANT, APT#, and COMPLEX, then an instance of S is as in Figure 5-1b. As the reader will notice, the schema S bears the same information that the schema R does, and thus by definition 5.1 schemas R and S are equivalent. This is in fact because the projection operation has been performed based on a data dependency, that is, the FD COMPLEX \rightarrow TYPE (corollary 3.1). The following example indicates a situation in which the projections of a relation carry more information (in fact, nonsense information) than the original relation does. Example 5.2: Consider the schema R of Example 5.1 and its snap-shot given in Figure 5-la. Now, if a schema S' consists of the projections of R on TYPE and APT#; and on TYPE, TENANT, and COMPLEX, then an instance of S' is as in Figure 5-2. As we see in Figure 5-2 each TENANT is associated with not only his APT but also with some other APTs of some other TENANTs. That is, the schema S' does not represent the same information as schema R, and thus assuming Definition 5.1 for representation, they are not equivalent.

Considering Definition 5.1, the FNF Algorithm always constructs a schema which satisfies the representation principle. This claim is formalized in the following theorem. <u>Theorem 5.1</u>. Let R be a schema constructed from a set of data dependencies G, using the FNF Algorithm, then R represents G.

An example of a relation and its projections

a) An instance of the schema R.

TACT (TENANT,	APT#,	COMPLEX,	TYPE)
Jack	1 7D	Nieman	Fur.
Mary	21 2H	Parkview	Unfur.
Phil	305F	Kraettli	Fur.
Mark	302A	Kraettli	Fur.
Mike	208C	Kraettli	Fur.
John	126 H	Const.	Unfur.
Tom	113 F	Const.	Unfur.

b) An instance of the schema S (projections of R).

CT (COMPLEX,	TYPE)	TAC (TENANT,	APT#,	COMPLEX)
Neiman	Fur.	Jack	17D	Nieman
Parkview	Unfur.	Mary	212 H	Parkview
Kraettli	Fur.	Phil	305F	Kraettli
Const.	Unfur.	Mark	302A	Kraettli
		Mike	208C	Kraettli
		John	12 6H	Const.
		Tom	11 3F	Const.

Proof: PART-1 of the algorithm constructs one relation scheme for each FD and one relation scheme for each MVD. Therefore. the schema constructed in PART-1 represents G. In PART-2, only those schemes that represent the same information as some other schemes of the schema are deleted, and thus the schema of PART-2 represents the schema of PART-1 which in turn represents G. PART-3 and PART-4 of the algorithm remove only extraneous attributes (i.e., attributes represented elsewhere in the schema, in stronger form) from relation schemes; and thus no information is lost in these two PARTs, and the schema constructed here (at the end of PART-4) represents the schema The other three parts of the algorithm, that is, of PART-2. PART-5, PART-6, and PART-7 decompose all of the decomposable schemes of the schema. All decomposition steps in these PARTs are based on the data dependencies holding in relation schemes, and thus considering Corollary 3.1, and assuming Definition 5.1, the decomposed schema (i.e., the final schema of the algorithm) represents the schema of PART-4 which in turn represents G, completing the proof. #

<u>Definition 5.2</u>. Let R and S be two relational schemas, then S represents the same information as R if it contains the same attributes and the same data dependencies as R.

Considering Definition 5.2, the FNF Algorithm does not always construct an schema satisfying the representation principle. In fact, the schema constructed by the first four PARTs of the algorithm satisfies the representation principle

An instance of the schema S'. (different projections of R of Figure 5-la).						
AT ($\underline{APT\#}$,	TYPE)	TTC (TENANT,	TYPE,	COMPLEX)		
17 D	Fur.	Jack	Fur.	Nieman		
212H	Unfur.	Mary	Unfur.	Parkview		
305F	Fur.	Phil	Fur.	Kraettli		
302A	Fur.	Mark	Fur.	Kraettli		
208C	Fur.	Mike	Fur.	Kraettli		
1 26H	Unfur.	John	Unfur.	Const.		
11 3F	Unfur.	Tom	Unfur.	Const.		

Figure 5-2

.

۰.

(Theorem 3.5), but PART-5, PART-6, and PART-7 may decompose relation schemes into subrelations violating some of the data dependencies. The following example clarifies this concept. Example 5.3: Given the following data dependencies G,

$$\begin{array}{ll} \mathbf{f}_{1}^{*} & \mathbf{A}, \mathbf{B} \longrightarrow \mathbf{C}, \mathbf{D}, \mathbf{E} \\ \mathbf{m}_{1}^{*} & \mathbf{C} \longrightarrow \mathbf{B} \end{array}$$

the algorithm first constructs the following schemes,

$$s_{1} \begin{cases} R_{1}(\underline{A}, \underline{B}, C, D, E) \\ R_{2}(\underline{C}, \underline{B}) \end{cases}$$

then it decomposes Rl into $R_{11}(\underline{C,B})$ and $R_{12}(\underline{A,C,D,E})$, and removes Rll because it is the same as R_2 .

$$s_{2} \begin{cases} R_{12} (\underline{A, C, D, E}) \\ R_{2} (\underline{C, B}) \end{cases}$$

Notice that the schema S_1 represents the same information as G (assuming either concepts -- Def. 5.1, or Def. 5.2), and the schema S_2 represents S_1 only if Definition 5.1 is considered. In other words, the databases of S_2 contain the same data as the databases of S_1 (satisfying Def. 5.1), but the FD A,B \rightarrow C,D,E is no longer represented by the schema S_2 (violating Def. 5.2).

One may suggest that instead of converting the schema S_1 into S_2 , we could simply remove the relation scheme R2 from S_1 , and thus, the schema S_1 consisting of R_1 satisfies both definitions of representation. The problem is, that as we discussed in Chapter II, the representation principle is not

the only principle that has to be satisfied by a "good" schema. The above suggestion, in fact, ignores the separation principle, and thus allows the schema to have undesirable properties.

One possible solution to the problem is to associate the dependency which is lost by a decomposition process to corresponding subrelations, by means of storing it in a specified area. This area then can be looked over when those subrelations need to be updated. Notice that the proposed solution can no longer be expedient if the above situation (i.e., losing a data dependency in a decomposition process) happens frequently in the schema design.

5.2.2 Separation Principle

Another goal in designing a "good" schema is to construct a schema in which the basic "units of information" are represented separately from each other. This is in fact, the intuitive motivation behind the normalization process [13]. In other words, the intention is to remove the undesirable properties (i.e., those properties that cause update anomalies discussed in Chapter I) from the schema by reducing the information to its more basic units. In this section, we will show that the schema constructed by the FNF Algorithm, is free from undesirable properties, and indeed it is in 4NF. This claim is formally illustrated by the following propositions and theorem.

<u>Proposition 5.1</u>. Considering a 2NF relation scheme $R(\underline{P}, N)$ where P is the primary key of R, and N is a non-empty set, then the MVD m: $X \rightarrow Y$ (or FD f: $X \rightarrow Y$) does <u>not</u> hold in R if both X and Y are subsets of P. <u>Proof</u>: (By contradiction) Since X and YeP, then relation R can be represented as $R(\underline{X}, \underline{Y}, \underline{A}, N)$, that is, we have $f_{\underline{P}}$: X,Y,A $\rightarrow N$

from f, and Axiom FD-MVD1 we get

m: $X, Y, A \longrightarrow N$

and if m: $X \rightarrow X$ holds in R, then from m and m₁ and Axiom MV04 we get

 $m_2: X, A \longrightarrow N$

from m_3 and f_1 and Axiom FD-MVD2 we have

 $f_2: X, A \rightarrow N$

which means N is partially dependent on key. This partial dependency is a contrary to the fact that the relation R is in 2NF.

<u>Proposition 5.2</u>. Given a 3NF relation scheme $R(\underline{P}, N)$ where p is the primary key of the relation R, and N is a non-empty set, then the MVD m: $X \longrightarrow Y$ (or FD f: $X \longrightarrow Y$) does not hold in R if both X and Y are subsets of N.

Proof: (By contradiction)

By our assumption relation R can be indicated as $R(\underline{P}, X, Y, B)$ from which we can conclude

 $f_1: P \rightarrow \{X, Y, B\}$

from f and Axiom FD6 we get

$$f; P \to Y$$

from f and m (assuming it <u>does</u> hold in R), and Axiom FD-MVD2 we get

 $f_2: X \longrightarrow Y$

which means attribute Y is transitively dependent on key, and thus violating 3NF schema characteristics, and contradiction to our assumption. #

<u>Proposition 5.3</u>. If relation scheme $R(\underline{P}, N)$ where P is the primary key of R and N is a non-empty set, is in 2NF, and if XeP and YeN, then MVD m: $X \longrightarrow Y$ (or FD f₁ X \longrightarrow) does not hold in R. <u>Proof</u>: (By contradiction)

If $P = \{X, A\}$ and $N = \{Y, B\}$ then we have R(X, A, Y, B) and we get

 $f_1: X, A \longrightarrow \{Y, B\}$

from f₁ and Axiom FD6 we get

 $f_2: X, A \longrightarrow Y$

from m (assuming it does hold in R) and f_2 , and Axiom FD-MVD2 we get

$$f_3: X \longrightarrow Y$$

which means a non-prime attribute is partially dependent on key, and thus violating the fact that R is in 2NF. # <u>Theorem 5.2</u>. Let R be a relational schema constructed using the FNF Algorithm, then R is in 4NF.

<u>Proof</u>: From Theorem 3.9 it follows that the schema constructed by the first five PARTs of the algorithm is at least in 3NF. PART-6 and PART-7 of the algorithm decompose any relation scheme in which a nontrivial multivalued dependency such as $X \longrightarrow Y$ holds, but not the functional dependency $X \longrightarrow C$ for every attribute C of that relation. This process which guarantees the schema to be in 4NF can be formalized as follows:

First, consider the set of relation schemes R' (that is, those relation schemes in which not <u>all</u> of the attributes make the key), and suppose $R_i(\underline{P}, N) \in R'$ where P is a set of attributes making the key of R_i . Now, assuming the MVD X \longrightarrow Y holds for R_i , the following possibilities should be considered:

- (1) X and Y ε P
- (2) X and Y ε N
- (3) $X \in P$ and $Y \in N$
- (4) $X \in N$ and $Y \in P$

Possibilities (1), (2), and (3) are not acceptable according to propositions (5.1), (5.2), and (5.3) respectively. For the case (4), relation R_i can be indicated as $R_i(\underline{Y},\underline{A},X,B)$. Since MVD X \longrightarrow Y is an element of the set G^{III} and $X \subseteq \Delta$ (i.e., X is a non-prime attribute), and $Y \subseteq K_{R_i}$ (i.e., a subset of key), then PART-6 of the algorithm decomposes R_i into R_i and R_i based on MVD X \longrightarrow Y. Now, MVD X \longrightarrow Y holds in R_i , and it is no longer a non-trivial MVD (because X U Y = U_{R_i1}).

By the above discussion we conclude that the set of relation schemes R' constructed by the first five PARTs of the algorithm is either already in 4NF, or it is converted to 4NF schema by PART-6.

Second, any relation scheme $R_j \in R^*$ (i.e., those relation schemes constructed merely from MVDs, and thus in each of them the entire set of attributes forms the key) is decomposed into its projections by PART-7 of the algorithm if a nontrivial MVD holds in it. Therefore, remaining relation schemes of R" are in 4NF, and thus the relational schema R which is R' U R" is in 4NF. #

5.2.3. Minimal Redundancy Principle

By the minimal redundancy principle (sometimes called nonredundancy principle) it is meant that although the constructed schema must represent all of the information of interest, it should not contain any redundant information (we will discuss the concept of redundancy later in this section).

As for the representation principle, different researchers have defined minimal redundancy in different ways. Here in this section, we will present and discuss our own definition of minimality which embodies a quite different view from those given by others. As a matter of fact, in other definitions, it is the relation scheme upon which attention is focused, to check whether it is required for the schema, or it is redundant. On the other hand, in our definition of minimality, we place stress on attributes rather than schemes. In fact, our intention is to avoid the repetition

of attributes in the schema as much as possible, and we call the schema with this characteristic <u>optimal</u>. The interesting point is, that the optimal schema can be achieved by our approach. We will return to this point after the following formal discussions of minimality. We begin with a standard definition given in literature, and then present our own definition of minimality.

<u>Definition 5.3</u> [6] Let R be a schema, then R is <u>minimal</u> if it does not contain any relation scheme R_i whose data dependencies are represented by the other schemas in R. <u>Definition 5.4</u> Let R be a schema, then R is <u>minimal</u> if for any other schema such as S which represents the same information as R, we have

 $\Sigma |R_i| \leq \Sigma |S_i|$

for all $R_i \in R$ for all $S_i \in S$ Example 5.4: Given the following schema S: $S\{R(\underline{A},\underline{B},\underline{C},\underline{D},\underline{E},\underline{F},\underline{G})\}$

and suppose two nontrivial MVDs $m_1: E \longrightarrow G$ and $m_2: C, D \longrightarrow F$ hold in R, then R can be decomposed into either $R_1(\underline{E},\underline{G})$ and $R_2(\underline{A},\underline{B},\underline{C},\underline{D},\underline{E},\underline{F})$, or $R_1'(\underline{C},\underline{D},\underline{F})$ and $R_2'(\underline{A},\underline{B},\underline{C},\underline{D},\underline{E},\underline{G})$ depending on the dependency chosen $(m_1 \text{ or } m_2)$ for decomposition. Therefore, the final schema will be either S' or S" (below).

$$\begin{cases} R_1 (\underline{E}, \underline{G}) \\ R_2 (\underline{A}, \underline{B}, \underline{C}, \underline{D}, \underline{E}, \underline{F}) \\ \\ R_1' (\underline{C}, \underline{D}, \underline{F}) \\ R_2' (\underline{A}, \underline{B}, \underline{C}, \underline{D}, \underline{E}, \underline{G}) \end{cases}$$

Notice that while both S' and S" represent the same information

as S, but only S' is the minimal schema. In fact, two attributes C and D are repeated in schema S", while in schema S' only one attribute E is repeated twice, and thus we have

$$\Sigma |S_{1}'| = |R_{1}| + |R_{2}| = 6 + 2 = 8$$

for all $S_{1}' \in S'$
$$\Sigma |S_{1}''|$$

for all $S_{1}'' \in S'' = |R_{1}'| = |R_{2}'| = 6 + 3 = 9$

In the above example, $\Sigma |S_i|$ is less than $\Sigma |S_i|$, because we decomposed S using MVD m₁ to get the schema S', and using MVD m₂ to get the schema S"; and this happened because $|\text{LEFT}_{m_1}| < |\text{LEFT}_{m_2}|$. This concept is formalized in the following proposition.

<u>Proposition 5.4</u>. Let R be the only decomposable relation scheme (i.e., some nontrivial data dependencies hold in it) in a schema S, and suppose there are alternatives for decomposing R (i.e., there are more than one nontrivial dependency holding in R), then decomposition leads to an optimal schema if the data dependency with the smallest number of attributes on its left side is chosen for decomposition. <u>Proof</u>: Suppose one of the dependencies holding in $R(\Delta)$ is m: $X_j \longrightarrow Y_j$, and suppose we decide to decompose R based on m, then we have

S'
$$\begin{cases} R_1(X_j, Y_j) \\ R_2(X_j, \Gamma) \text{ where } = \Gamma - \Delta X_j - Y_j \end{cases}$$

and thus we get

$$\Sigma |S_{i}^{i}| = |R_{1}| + |R_{2}| = |X_{j}| + |Y_{j}| + |X_{j}| + |X_{j}|$$

and since $|\Delta|$ is constant, then the value of $\Sigma |S_i|$ depends only on $|X_j|$, and thus it is minimum if $|X_j|$ is minimum, completing the proof. #

Notice that Proposition 5.4 assumes a schema with only one decomposable scheme. Indeed, if there are more than one decomposable schemes in a schema, then choosing a data dependency with the smallest left side does not necessarily lead to an optimal schema. The following example clarifies this concept.

Example 5.5: Given the following data dependencies:

 $m_{1}: A, B, D, E, F \longrightarrow C, G$ $m_{2}: D, E, F \longrightarrow I, J$ $m_{3}: D, E \longrightarrow F$ $m_{4}: J \longrightarrow E$

we can first construct the following schema:

 $R_{1}(\underline{A},\underline{B},\underline{D},\underline{E},\underline{F},\underline{C},\underline{G})$ $R_{2}(\underline{D},\underline{E},\underline{F},\underline{I},\underline{J})$

$$R_{3}(\underline{D},\underline{E},\underline{F})$$
$$R_{4}(\underline{J},\underline{E})$$

Here, relation schemes R_1 nad R_2 are decomposable; No alternatives exist for R_1 , and it can be decomposed only based on dependency m_3 . Thus, the projections are $R_{11}(\underline{D},\underline{E},\underline{F})$ and $R_{12}(A,B,D,E,C,G)$, and since R_{11} is the same as R_3 it is remoced, and we have

 $R_{12}(\underline{A},\underline{B},\underline{D},\underline{E},\underline{C},\underline{G})$ $R_{2}(\underline{D},\underline{E},\underline{F},\underline{I},\underline{J})$ $R_{3}(\underline{D},\underline{E},\underline{F})$ $R_{4}(\underline{J},\underline{E})$

For R_2 , there are two alternatives, it can be decomposed either based on dependency m_3 or dependency m_4 . If we decompose R_2 based on dependency m_4 (with only one attribute on the left) into $R_{21}(\underline{J},\underline{E})$ and $R_{22}(\underline{D},\underline{F},\underline{I},\underline{J})$, and remove R_{21} which is the same as R_4 , then we get

s'
$$\begin{cases} R_{12}(\underline{A}, \underline{B}, \underline{D}, \underline{E}, \underline{C}, \underline{G}) \\ R_{22}(\underline{D}, \underline{F}, \underline{I}, \underline{J}) \\ R_{3}(\underline{D}, \underline{E}, \underline{F}) \\ R_{4}(\underline{J}, \underline{E}) \end{cases}$$

But, if we decompose R_2 based on dependency m_3 (with two attributes on the left) into $R_{21}'(\underline{D},\underline{E},\underline{F})$ and $R_{22}'(\underline{D},\underline{E},\underline{I},\underline{J})$, and remove R_{21}' which is the same as R_3 , then we get

$$R_{12}(\underline{A},\underline{B},\underline{D},\underline{E},\underline{C},\underline{G})$$

$$R_{22}'(\underline{D},\underline{E},\underline{I},\underline{J})$$

$$R_{3}(\underline{D},\underline{E},\underline{F})$$

$$R_{4}(\underline{J},\underline{E})$$

and now, by decomposing R_{22} ' (no alternative) into R_{22} ' $l(\underline{J},\underline{E})$ and R_{22} ' $2(\underline{D},\underline{I},\underline{J})$, and then removing R_{22} 'l which is the same as R_4 we get

s"
$$\begin{cases} R_{12} (\underline{A}, \underline{B}, \underline{D}, \underline{E}, \underline{C}, \underline{G}) \\ R_{22} ' 2 (\underline{D}, \underline{I}, \underline{J}) \\ R_{3} (\underline{D}, \underline{E}, \underline{F}) \\ R_{4} (\underline{J}, \underline{E}) \end{cases}$$

Notice that, although we didn't choose the dependency with the smallest left side to get the schema S", but S" \underline{is} optimal. This is because

$$\Sigma |S_{i}^{"}| = 14$$
for all $S_{i}^{"} \in S^{"}$

$$\Sigma |S_{i}^{'}| = 15$$
for all $S_{i}^{'} \in S^{'}$

and thus

 $\Sigma |S_{i}^{"}|$ < $\Sigma |S_{i}^{"}|$ for all $S_{i}^{"} \in S$ for all $S_{i}^{'} \in S^{'}$

This happened because, in decomposition process of S', we used two different data dependencies to decompose R_1 and R_2 , while both relation schemes could be decomposed based on

the same data dependency (MVD m_2), and thus generating the same schemes that could be removed from the schema (in fact, it is what we did for constructing S"). This claim can be formally illustrated as follows:

Theorem 5.3. Let S be an schema containing some decomposable schemes, then decomposition leads to an optimal schema if:

- each decomposable scheme that has no dependency in common with other schemes, is decomposed based on the dependency with the smallest left side, and
- (2) all decomposable schemes that have some dependencies in common, are decomposed based on the common dependency with the smallest left side.

<u>Proof</u>: The proof for the part(1) of the theorem directly follows from Proposition 5.4, and part(2) can be proved as follows:

Suppose S is a schema, and suppose $R \subseteq S$ with |R| = p is a subschema consisting merely of the schemes that have some data dependencies holding in them in common, then using one of the common data dependencies, say, m: $X_{j} \longrightarrow Y_{j}$, we can decompose them as follows:

$$R \begin{cases} R_{1}(\Lambda_{1}) & \xrightarrow{X_{j} \longrightarrow Y_{j}} \\ \xrightarrow{X_{j} \longrightarrow Y_{j}} & \begin{cases} R_{11}(X_{j}, Y_{j}) \\ R_{12}(X_{j}, Y_{1}) & \text{where } \Gamma 1 = \Lambda_{1} - \\ & (X_{j} + Y_{j}) \end{cases} \\ R_{2}(\Lambda_{2}) & \xrightarrow{X_{j} \longrightarrow Y_{j}} & \begin{cases} R_{21}(X_{j}, Y_{j}) \\ R_{22}(X_{j}, 2) & \text{where } \Gamma 2 = \Lambda_{2} - \\ & (X_{j} + Y_{j}) \end{cases} \\ \vdots \\ R_{p}(\Lambda_{p}) & \xrightarrow{X_{j} \longrightarrow Y_{j}} & \begin{cases} R_{p1}(X_{j}, Y_{j}) \\ R_{p2}(X_{j}, \Gamma_{p}) & \text{where } \Gamma_{p} = \Lambda_{p} - \\ & (X_{j} + Y_{j}) \end{cases} \end{cases}$$

Now, since all relation schemes R_{il} (for $l \le i \le p$) are the same, we remove all but R_{ll} from the schema, and therefore we have

$$R' \begin{cases} R_{11} (X_{j}, Y_{j}) \\ R_{12} (X_{j}, \Gamma_{1}) \\ R_{22} (X_{j}, \Gamma_{2}) \\ R_{32} (X_{j}, \Gamma_{3}) \\ \vdots \\ Rp_{2} (X_{j}, \Gamma_{p}) \end{cases}$$

and hence, we get

$$\Sigma |\mathbf{R}_{i}^{t}| = |\mathbf{x}_{j}| + |\mathbf{Y}_{j}|$$
for all $\mathbf{R}_{i}^{t} \in \mathbf{R}^{t} + |\mathbf{x}_{j}| + |\Delta_{1}| - |\mathbf{x}_{j}| - |\mathbf{Y}_{j}|$

$$+ |\mathbf{x}_{j}| + |\Delta_{2}| - |\mathbf{x}_{j}| - |\mathbf{Y}_{j}|$$

$$+ |\mathbf{x}_{j}| + |\Delta_{3}| - |\mathbf{x}_{j}| - |\mathbf{Y}_{j}|$$

$$\vdots$$

$$+ |\mathbf{x}_{j}| + |\Delta_{p}| - |\mathbf{x}_{j}| - |\mathbf{Y}_{j}|$$

$$\Sigma |\mathbf{R}_{i}^{t}| = \epsilon \mathbf{R}^{t} = (|\Delta_{1} + \Delta_{2} + \Delta_{3} + \dots + \Delta_{p}|) + |\mathbf{X}_{j}| - (p - 1) \cdot |\mathbf{Y}_{j}|$$

$$= \sum_{m=1}^{p} |\Delta_{n}| + |\mathbf{x}_{j}| - (p - 1) \cdot |\mathbf{Y}_{j}|$$

$$(5.2.3.2)$$

On the other hand, if we decompose each relation scheme in R, using any arbitrary dependency not in common with other schemes, then we get

$$\mathbb{R}^{\mathsf{t}} \left\{ \begin{array}{c} \mathbb{R}_{1} \left(\Delta_{1} \right) \xrightarrow{\mathbb{A}_{1} \longrightarrow \mathbb{B}_{1}} \mathbb{R}_{1} \left(\mathbb{A}_{1}, \mathbb{B}_{1} \right) \\ \mathbb{R}_{12} \left(\mathbb{A}_{1}, \Gamma_{1}^{\mathsf{t}} \right) \text{ where } \Gamma_{1}^{\mathsf{t}} = \Delta_{1} - \left(\mathbb{A}_{1} + \mathbb{B}_{1} \right) \\ \mathbb{R}_{2} \left(\Delta_{2} \right) \xrightarrow{\mathbb{A}_{2} \longrightarrow \mathbb{B}_{2}} \mathbb{R}_{22} \left(\mathbb{A}_{2}, \mathbb{B}_{2} \right) \\ \mathbb{R}_{22} \left(\mathbb{A}_{2}, \Gamma_{2}^{\mathsf{t}} \right) \text{ where } \Gamma_{2}^{\mathsf{t}} = \Delta_{2} - \left(\mathbb{A}_{2} + \mathbb{B}_{2} \right) \\ \vdots \\ \mathbb{R}_{p} \left(\mathbb{A}_{p} \right) \xrightarrow{\mathbb{A}_{p} \longrightarrow \mathbb{B}_{p}} \left\{ \mathbb{R}_{p1} \left(\mathbb{A}_{p}, \mathbb{B}_{p} \right) \\ \mathbb{R}_{p2} \left(\mathbb{A}_{p}, \Gamma_{p}^{\mathsf{t}} \right) \text{ where } \Gamma_{p}^{\mathsf{t}} = \Delta_{p} - \left(\mathbb{A}_{p} + \mathbb{B}_{p} \right) \end{array} \right\}$$

and thus we get

$$\Sigma |R_{i}^{"}| = |A_{1}| + |B_{1}| + |A_{1}| + |\Delta_{1}| - |A_{1}| - |B_{1}|$$

for all $R_{i}^{"} \in R^{"}$
$$|A_{1}| - |B_{1}|$$
$$+ |A_{2}| + |B_{2}| + |A_{2}| + |\Delta_{2}| - |A_{2}| - |B_{2}|$$
$$+ |A_{3}| + |B_{3}| + |A_{3}| + |\Delta_{3}| - |A_{3}| - |B_{3}|$$
$$|A_{3}| - |B_{3}|$$
$$\vdots$$
$$\vdots$$
$$+ |A_{p}| + |B_{p}| + |A_{p}| - |\Delta_{p}| - |A_{p}| - |A_{p}| - |B_{p}|$$

$$\Sigma |\mathbf{R}_{i}^{"}| = (|\Delta_{1} + \Delta_{2} + \Delta_{3} + \dots + \Delta_{p}|) +$$

for all $\mathbf{R}_{i}^{"} \in \mathbf{R}^{"}$
$$|\mathbf{A}_{1}| + |\mathbf{A}_{2}| + \dots + |\mathbf{A}_{p}|$$

(5.2.3.3)

$$\Sigma |\mathbf{R}_{i}^{"}| = \sum_{m=1}^{p} |\Delta_{m}| + \sum_{m=1}^{p} |A_{m}| \qquad (5.2.3.4)$$

for all $\mathbf{R}_{i}^{"} \in \mathbf{R}^{"} = \sum_{m=1}^{m} |\Delta_{m}| = 1$

$$\Sigma |\mathbf{R}_{1}^{"}| = \sum_{m=1}^{p} |\Delta_{m} + A_{m}| \qquad (5.2.3.5)$$

for all $\mathbf{R}_{1}^{'} \in \mathbf{R}^{'}$

Now, by analyzing formulas (5.2.3.2) and (5.2.3.4), we can conclude that choosing a common dependency (if any) for decomposition is always better (in the sense of optimality) than considering an arbitrary dependency, even if it has the smallest left side in the set of data dependencies. Indeed, we need to prove that for all situations, the value of (5.2.3.2) is less than the value of (5.2.3.4).

Consider the formula (5.2.3.2) which calculates the number of attributes (including repetitions) of a schema (or a subschema), for a situation that all of the relation schemes of the schema are decomposed based on the same data dependency (i.e., a dependency holding in all of them). Now, assume that one of these schemes is not decomposed based on the common dependency, but it is decomposed based on an arbitrary dependency. Then, as it is shown below, there will be an increase in the number of attributes of the schema.

(a) if all
$$R_j \in R$$
 are decomposed based on the same
data dependency $X_j \longrightarrow Y_j$:

$$\sum_{j} |R_{j}| = \sum_{m=1}^{\infty} |\Delta_{m}| + |X_{j}| - (p-1) \cdot |Y_{j}|$$

for all $R_{j} \in R$ m=1 (a.1)

(b) if all $R_i \in R$ (except R_k) are decomposed based on the same data dependency $X_j \longrightarrow Y_j$, and R_k is decomposed based on an arbitrary dependency $A_k \longrightarrow B_k$:

$$\begin{split} \Sigma |R_{i}| & p_{i} |\Delta m| + |X_{j}| - I(p-1) - 1] \\ \text{for all } R_{i} \in R - R_{k} & m=1, \ m\neq k \\ \cdot |Y_{j}| & (b.1) \\ |R_{k}| = |A_{k}| + |B_{k}| + |A_{k}| + |\Delta_{k}| - (|A_{k}| + |B_{k}|) = \\ & |A_{k}| + |\Delta_{k}| & (b.2) \end{split}$$

if we add (b.1) and (b.2) together we get

$$\Sigma |R_{i}| = \sum_{m=1}^{p} |\Delta_{m}| + |X_{j}| - (p - 2) \cdot |Y_{j}| + |A_{k}| = \sum_{m=1}^{p} |\Delta_{m}| + |\Delta_{k}| = |A_{k}| + |\Delta_{k}| + |A_{k}| + |A_{k}|$$

and we conclude that, (b.3) is always greater than (a.1), by indicating that the expression [(b.3) - (a.1)] is always positive.

$$\sum |R_i| \text{ of (b.3)} - \sum |R_i| \text{ of (a.1)} = |A_k| + |Y_j|$$

for all $R_i \in R$ for all $R_i \in R$

As conclusion, the number of attributes increases at least by 2 (when $|A_k| = |Y_j| = 1$). By the same token, if the number of relation schemes that are decomposed based on arbitrary dependencies is p', then the number of attributes in the schema increases by

$$\sum_{k=1}^{p'} |A_k| + p' \cdot |Y_j|.$$

The above analysis and the Proposition 5.4, lead to the conclusion that, when there are choices for decomposition, we can construct an optimal schema if: (i) those schemes that have some data dependencies in common, are decomposed based on the common dependency with the smallest left side, completing the proof. #

Considering Proposition 5.4 and Theorem 5.3, now it is possible to modify the FNF Algorithm to construct an optimal schema. To do so, we actually need to modify only those PARTs of the algorithm in which decompositions take place. In the FNF Algorithm, a decomposable scheme is decomposed based on a data dependency holding in it which is encountered first by the algorithm. Now, since the intention is the construction of an optimal schema, the algorithm should categorize all decomposable schemes before practically decomposing any scheme. Then considering Theorem 5.3, it should decompose each relation scheme based on an appropriate data dependency. One suggestion which leads to an efficient implementation of the problem is given in the following algorithm.

5.2.3.1. Description of the Optimal Decomposition Algorithm

For this (Figure 5-3) algorithm we need to construct a m x n matrix M, where m is the number of decomposable schemes R in the schema and n is the cardinality of the set of data dependencies G which is to be checked for decomposition. Then, $M_{i,j} = 1$, if the data dependency $g_j \in G$ holds for relation scheme $R_i \in R$ (i.e., if R_i can be decomposed based on g_j) and $M_{i,j} = 0$ otherwise. In addition we need to define an array A of n elements, so that A_k states the number of decomposable schemes that data dependency g_k holds for them.

Algorithm 5-1

An Algorithm to Decompose a Set of Decomposable Schemes

<u>into an Optimal Form</u>

- INPUT: A set R of m decomposable relation schemes; and a set G of n data dependencies.
- OUTPUT: A set of decomposed schemes in optimal form.
- STEP 1: Find the element of array A, say A_k, which contains the largest number (if more than one element.found, then choose the one such that the corresponding dependency g_k has the smallest left side), then remove this element.
- STEP 2: Decompose all relation schemes in which data dependency g_k holds. These schemes can be found simply by scanning column k of matrix M. That is, any row that has a 'l' in column k, corresponds to one of these schemes. Mark these schemes to indicate that they have been decomposed. Then update array A, that is, for any R_i ε R that was just marked and M_{i,j} = 1, decrease A_j by one.
- STEP 3: Repeat STEP 1 and STEP 2 until all R_i ε R are marked, that is, all decomposable schemes are decomposed.

Figure 5-3

5.3 Bernstein's Synthetic Approach

There have been several approaches for designing an algorithm to synthesize the relational schema from a set of functional relationships. Wang and Wedekind [29] proposed such an algorithm to produce third normal form relations from a given set of functional dependencies. However, their algorithm could generate two relations with keys that are functionally equivalent. Another approach has been made by Bernstein [8]. Bernstein's work was based on the approach given by Delobel and Casey [15]. The relational schema produced by his algorithm is in third normal form and contains a minimal number of rela-The fourth normal form schema cannot be achieved tions. by this algorithm because, only functional dependencies have been considered and multivalued dependencies have been ignored. This algorithm is indicated in Figure 5-4. Example 5-6: Given the following functional dependencies:

$$f_{1}: A \rightarrow X$$

$$f_{2}: A, B \rightarrow Y$$

$$f_{3}: A, B \rightarrow Z$$

$$f_{4}: B, Z \rightarrow A$$

$$f_{5}: B, Z \rightarrow Y$$

$$f_{6}: B, Y \rightarrow A$$

$$f_{7}: B, Y \rightarrow Z$$

The nonredundant covering of the above set of FDs is as follows:

$$f_1: A \rightarrow X$$
Algorithm 5-2

Bernstein's Algorithm to synthesize 3NF schema

- INPUT : A set F of FDs.
- OUTPUT: A set of relation schemes in 3NF.
- STEP 1: Find a nonredundant covering for F. Call this set G.
- STEP 2: Partition all functional dependencies in G, into groups where all of the FDs in each group have identical left sides.
- STEP 3: Let G_1 and G_2 be any pair of groups, then merge these two groups together if there exists a bijection $X \leftrightarrow Y$ (X and Y are left sides of groups G_1 and G_2 respectively) in the closure of G, G⁺.
- STEP 4: Construct a relation scheme for each group by taking the set union of all the attributes appearing in that group. The set of attributes that appears on the left side of any FDs of that group is the key of the relation scheme.

Figure 5-4

$$f_{2}: A, B \longrightarrow Y$$
$$f_{4}: B, Z \longrightarrow A$$
$$f_{7}: B, Y \longrightarrow Z$$

Notice that f_3 is redundant because it could be derived from f_2 and f_7 using Axiom FD4. Similarly f_5 and f_6 are redundant because f_5 is implied by f_4 and f_2 , and f_6 is implied by f_7 and f_4 .

Now, from the above set of nonredundant covering, the algorithm constructs the following third normal form relational schema:

 $R_{1}(\underline{A}, X)$ $R_{2}(\underline{A}, \underline{B}, Y)$ $R_{3}(\underline{B}, \underline{Z}, A)$ $R_{4}(\underline{B}, \underline{Y}, Z)$

An important problem concerning the Algorithm 6-1 is finding the nonredundant covering of the set of FDs (STEP 1) efficiently. Recall from Chapter II that an FD $f_i \in F$ is said to be redundant in F if $F^+ = (F - \{f_i\})^+$. In other words, if $f_i \in (F-\{f_i\})^+$, then f_i is redundant and can be removed from the set F without altering the closure of the set. To find the nonredundant covering of F, we need to find the FDs that are redundant, and then remove them from the set. Therefore, an algorithm which can decide whether or not a single functional dependency is in the closure of a given set of FDs, appears to be essential for synthesizing relations from FDs. A membership algorithm is given in the appendix.

5.4 Fagin's Decomposition Approach

In [18], Fagin has presented a decomposition approach for constructing the relational schema. In this approach, he considers both kind of dependencies, functional and multivalued, and the schema constructed is in 4NF. He begins his normalization process by forming a single relation scheme consisting of all attributes of the data base. Then this relation (if not in 4NF) would be decomposed into two of its projections. This process continues for each subrelation not in 4NF, until all of them are in 4NF. The following example indicates a 4NF normalization process.

Example 5.7: Reconsider the data base of the Example 3.12 given in Chapter III. The data dependencies are as follows:

fl:	CLASS, SECTION \rightarrow INSTRUCTOR			
f ₂ :	CLASS, SECTION, DAY \rightarrow ROOM			
f ₃ :	STUDENT \rightarrow MAJOR, YEAR			
f ₄ :	INSTRUCTOR \rightarrow RANK, SALARY			
^m 1:	CLASS, SECTION $\rightarrow \rightarrow$ STUDENT, MAJOR, EXAM, YEAR			
^m 2 [:]	CLASS, SECTION> INSTRUCTOR, RANK, SALARY			
^m 3:	CLASS, SECTION $\rightarrow \rightarrow$ DAY, ROOM			
^m 4:	CLASS \rightarrow TEXT			
m ₅ :	CLASS, SECTION, STUDENT \rightarrow EXAM			

The normalization process begins by forming a single relation scheme R containing all attributes:

R(CLASS, SECTION, STUDENT, MAJOR, EXAM, YEAR, INSTRUCTOR,

RANK, SALARY, TEXT, DAY, ROOM)

Based on the MVD m_1 the relation scheme R can be decomposed into R_1 and R_2 as follows:

 R_1 (CLASS, SECTION, STUDENT, MAJOR, EXAM, YEAR)

 R_2 (CLASS, SECTION, INSTRUCTOR, RANK, SALARY, TEXT, DAY,

ROOM)

Neither R_1 nor R_2 is in 4NF. R_1 can be decomposed based on m_5 into

R₁₁ (CLASS, SECTION, STUDENT, EXAM)

R12 (CLASS, SECTION, STUDENT, MAJOR, YEAR)

Although R_{11} is in 4NF, R_{12} is not. It can be decomposed based on FD f_3 into

R₁₂₁ (STUDENT, MJAOR, YEAR)

R₁₂₂ (CLASS, SECTION, STUDENT)

Both R_{121} and R_{122} are in 4NF. We now decompose R_2 based on MVD m_2 into

R₂₁ (CLASS, SECTION, INSTRUCTOR, RANK, SALARY)

R22 (CLASS, SECTION, TEXT, DAY, ROOM)

Considering FD f_4 , R_{21} can be decomposed into

R₂₁₁ (INSTRUCTOR, RANK, SALARY)

R₂₁₂ (CLASS, SECTION, INSTRUCTOR)

And using MVD m_4 , R_{22} can be decomposed into

R₂₂₁ (CLASS, TEXT)

R₂₂₂ (CLASS, SECTION, DAY, ROOM)

Now, the schema consisting of the set of relation schemes $\{R_1, R_{121}, R_{122}, R_{211}, R_{212}, R_{221}, R_{22}\}$ is in 4NF. If we

remove R₁₂₂ from the schema (because it is a projection
of R₁₁, and therefore it is redundant) we get the same
set of relation schemes that we got earlier in Chapter III,
using the new method.

As we saw in Section 5.3 and in this section, the Bernstein's synthetic approach guarantees 3NF schema and not higher, while the decomposition process can lead to a 4NF family. In addition, synthesized schema satisfies the representation principle of Definition 5.2 given in Section 5.2, and decomposition leads to a schema which satisfies the representation principle defined by Definition 5.1 [6]. Although the schema constructed by Fagin's decomposition method is in 4NF, it is not necessarily optimal. On the other hand, as it was shown in previous chapters, the dependency-to-relation approach not only guarantees the 4NF schema, it also can lead to an optimal schema. In addition, this method does not violate the representation principle of Definition 5.1, and avoids violating the Definition 5.2 as much as possible.

Considering the above discussion and the idea given in [6], we can summarize the differences in the following table:

Method	Bernstein's	Fagin's	Authors
Normal Form	3NF .	4NF	4NF
Data Dependencies	FDs	FDs+MVDs	FDs + MVDs
Definition of Minimality	Def.5.3		Both, Def. 5.3 and Def. 5.4
Definition of Representation	Def. 5.2	Def.5.1	Def. 5.1 and mostly Def. 5.2

5.5. Chapter Summary and Remarks

Properties of the relation schemes constructed by the FNF Algorithm have been discussed in detail. It has been proven that the schema constructed is in 4NF and thus free of undesirable properties. It has been also shown that the schema may possibly contain some redundant information, and thus it is not optimal. To solve this problem, heuristics have been employed and "optimal" has been defined. Furthermore, an algorithm (which can be efficiently implemented) has been presented for decomposition of a set of relation schemes into an optimal form. Finally, properties of relations by the FNF Algorithm and by the other approaches have been compared.

CHAPTER VI

SUMMARY AND CONCLUSION

6.1 Summary

It was the purpose of this thesis to investigate the basic concepts of the relational model, and to present an approach for constructing relational schema from functional and multivalued dependencies. Although some of the previous works [8,11,14,18, et al.] have addressed a similar problem, these approaches present a number of situations that are not desirable for the data base. Many of these problems which will be pointed out in the next section have been resolved by the new method presented in this thesis.

The theoretical background for the approach was provided in Chapter II. The algebraic rules for both kind of dependencies (functional and multivalued), were examined. We also proposed two special closures for mixed dependencies (i.e., FDs and MVDs), and demonstrated their importance for the relational schema design. Two algorithms for computation of these closures along with the extensive analyses of them were given in Chapter IV.

The main algorithm of the approach was presented in

Chapter III. The process was started with designing a simple algorithm to construct one relation scheme for each explicitly designated dependency. We then examined this simple algorithm, and discovered a number of problems presented by it. These problems which were imputed to the ignorance of the composing rules for data dependencies, led us to modify the algorithm. Finally, after an iterative refinement of the algorithm, we presented the main algorithm (Fourth Normal Form Algorithm).

In Chapter V we examined the important properties of the schema constructed by the FNF Algorithm.

We proposed (in Section 5.2) a new concept for optimality, and made an extensive investigation toward constructing optimal relational schema.

Finally, in sections 5.3 and 5.4, we examined the Bernstein's synthetic algorithm and Fagin's decomposition method. We also compared their methods with the new approach presented in this dissertation.

6.2 Analysis of the Approach

First of all, since both kind of dependencies have been considered for the design approach, the constructed schema is in 4NF, and therefore many of the problems concerning the synthesized 3NF schema have been eliminated.

In addition, allowing the data base administrator

to inspect the schema at any desirable intermediate level, leads to practical results as well as theoretical consequences. This intervention can be performed efficiently since it might be categorized based on the different parts of the FNF algorithm. In other words, the result of each part of the algorithm can be examined independently -one part eliminates explicit partial dependencies, another part concentrates on implicit partial dependencies, another part is concerned with transitive dependencies, etc.. As an example, consider a situation in which the set of dependencies designated by the data base administrator includes the following:

- $f_3: EMPLOYEE \longrightarrow SALARY.$

Syntactically, the dependency f_3 can be derived from dependencies f_1 and f_2 using Axiom FD3. Note that what actually is implied by f_1 and f_2 is an FD in which EMPLOYEE determines the SALARY of the MANAGER, and this is completely different from f_3 in which EMPLOYEE determines its own SALARY. Therefore, f_3 is not redundant and should not be removed from the schema. This problem can be overcome by considering a signal in PART-8 when implementing the FNF Algorithm. This signal warns the data base administrator about deleting a nonredundant relation from the schema, and hence he can make the proper decision regarding this matter.

Moreover, using the concepts given in Subsection 5.2.3, the FNF Algorithm can be further modified to produce an optimal schema. This could save a large amount of storage for storing the data, because the repetition of the attributes has been minimized.

Finally, since the normalization process is performed in a step by step manner (i.e., generating lNF schema, then transforming the lNF into 2NF, 2NF into 3NF, etc.), the results can be used for the analysis of different normal forms and their transformations.

6.3 Suggestions for Further Work

This dissertation dealt solely with the logical issue of the relational data base. Specifically speaking, we focused more on the logical significance of the proposed algorithms than on their efficiency. Therefore, an improvement to the approach will be desirable regarding this matter.

Although this work provides all relevant algorithms, no actual programming has been done. These algorithms, especially the FNF Algorithm and those that are used for constructing the closure II and closure III, should be tested on real data bases, and their performance evaluated.

Another direction worth investigating is improving the approach to construct the relational schema satisfying

the representation principle of Definition 5.1, and not violating the Definition 5.2.

Finally, this dissertation leaves the designating of functional and multivalued dependencies completely in the hands of the data base administrator. Some work is needed for the systematic detection of these dependencies.

BIBLIOGRAPHY

- Aho, A.V., C. Beeri, and J.D. Ullman, "The theory of joins in relational databases", ACM Transactions on database Systems, September 1979.
- Aho, A.V., J.E. Hopcroft and J.D. Ullman, "The Design and Analysis of Computer Algorithms", Addison Wesley, Mass., 1974.
- Armstrong, W.W., "Dependency structure of data base relationships", Proc. 74 IFIP, North-Holland, 1974.
- Beeri, C., "On the membership problem for multivalued dependencies in relational databases", TR229, Dept. of EECS, Princeton University, 1977.
- Beeri, C. and P.A. Bernstein, "Computational problems related to the design of normal farm relation schemes", ACM Transactions on Database Systems, 1979.
- Beeri, C., P.A. Bernstein, and N. Goodman, "A sophisticate's introduction to database normalization theory", Proc. ACM Intl. Conf. on very large Data Bases, 1978.
- Beeri, C., R. Fagin, and J.H. Howard, "A complete axiomatization for functional and multivalued dependencies", ACM/SIGMOD International Symposium on Management of Data, 1977.
- Bernstein, P.A. "Synthesizing third normal form relations from functional dependencies", ACM Trans. on Database Systems. 1976.
- Cadiou, J.M. "On Semantic Issues in the Relational Model of Data", IBM Research Laboratory, San Jose, CA, 1977.
- 10. Codd, E.F. "A relational model for large shared data banks", CACM, 1970.

- 11. Codd, E.F. "Further normalizations of the data base relational model", In Data Base Systems, Courant Inst. Computer Science Symposium 6, R. Rustin, Ed., Prentice-Hall, Englewood Cliffs, NJ, 1972.
- 12. Codd, E.F. "Relational Completeness of Data Base Sublanguages", Courant Computer Science Symposia Series, Vol. 6, Englewood Cliffs, NJ, Prentice-Hall, 1972.
- Data, C.J. "An Introduction to Database Systems, Addison Wesley, Mass., 1977.
- 14. Delobel, C., "Normalization and hierarchical dependencies in the relational data model", ACM Transactions on Database Systems, 1978.
- 15. Delobel, C. and R.C. Casey, "Decomposition of a database and the theory of Boolean Switching functions", IBM J. Res., 1972.
- 16. Fagin, R. "Functional dependencies in a relational database and propositional logic", IBM Res. Laboratory, San Jose, CA, 1976.
- Fagin, R. "The decomposition versus synthetic approach to relational database design. Proc. Third Intl. Conf. on very large data bases, Tokyo, October 1977.
- 18. Fagin, R., Multivalued dependencies and a new norman form for relational databases", ACM Trans. on Database Systems, 1977.
- 19. Fagin, R., "Functional dependencies in a relational database and propositional logic", IBM J. Res. and Develop., November 1977.
- Fadious, R., "Mathematical foundations for relational data bases", doctoral diss., Michigan State Univ., 1975.
- 21. Hutt, A.T.F., "A relational data base management system", John Wiley & Sons, NY 1979.
- 22. Hagihard, K., M. Ito, K. Tanighuchi and T. Kasami, "Decision problems for multivalued dependencies in relational databases", SIAM J. Comput., May 1979.
- 23. Luk, W.S., "Possible membership of a multivalued dependency in a relational database", Information Processing letters, August 1979.

- 24. Meldelzon, A.O. "On axiomatizing multivalued dependencies in relational databases" J. ACM, 1979.
- 25. Martin, J., "Computer data-base organization", Prentice-Hall, Englewood Cliffs, NJ 1975.
- 26. Rissanen, J. "Independent components of relations", ACM Trans. on Database Systems, 1977.
- 27. Schmid, H.A. and J.R. Swenson, "On the semantics of the relational model", ACM/SIGMOD International Symposium on Management of Data, 1976.
- 28. Stant, D.F., and D.F. McAllister, "Discrete Mathematics in Computer Science", Prentice-Hall, 1977.
- 29. Wang, C.P. and H.H. Wedekind, "Segment Synthetics in Logical Data Base Design", IBM J. of Res. and Dev., January 1975.
- 30. Wiederhold, G. "Data base design", McGraw-Hill, NY, 1977.
- 31. Zaniolo, C. "Analysis and design of relational schemata for database systems", doctoral diss., UCLA, 1976.

APPENDIX

Completeness of the inference rules for FDs and MVDs

Using the definition it is easy to verify the validity of the inference rules. To prove the completeness of this set of rules we need to show that each dependency which is implied by $F \bigcup M$ (F is a set of FDs, and M is a set of MVDs) belongs to $(F,M)^+$.

Recall that a dependency f is implied by $F \bigcup M$ if there is no counterexample relation such that all dependencies in $F \bigcup M$ are valid in it but f is not. To show completeness of the rules, we have to show that for each dependency not in $(F,M)^+$ such a counterexample relation does exist. Eafore we present the completeness theorem we need a few concepts.

Let X be a subset of U. There are several sets Y such that the MVD X \longrightarrow Y is in (F,M)⁺, (e.g., X \longrightarrow U-X is always in (F,M)⁺. Following Fagin, we use the notation X \longrightarrow Y₁|Y₂|...|Y_k to denote the collection of MVD's X \longrightarrow Y₁, X \longrightarrow Y₂,..., X \longrightarrow Y_k. From now on,

[†]This discussion followsquite closely that of Beeri, Fagin and Howard [7].

when this notation is used, we assume that none of the sets Y_1, \ldots, Y_k is the empty set.

Let us denote by DEP(X) the family of all sets Y for which $X \rightarrow Y$. (DEP(X) is, of course, a function of the given sets of dependencies F and M.) We have seen that DEP(X) is closed under boolean operations (MVD5,MVD6). Therefore, it contains a unique subfamily with the following properties:

- a) The sets in the subfamily are nonempty.
- b) Each pair of sets in the subfamily is disjoint.
- c) Each set in DEP(X) is a union of sets from the subfamily.

This subfamily consists of all nonempty minimal sets in DEP(X), i.e., those sets that do not contain any other nonempty set of DEP(X). We call this subfamily the <u>depen-</u> <u>dency basis</u> of X. If Y_1, \ldots, Y_k are the sets in the dependency basis of X, then as Fagin noted [18], "the generalized" MVD X $\longrightarrow Y_1 | \ldots | Y_k$ contains all the information about MVDs that have X as their left side.

Let X* denote the set of all attributes that are functionally dependent on X (by functional dependencies in $(F,M)^+$). Clearly X \subseteq X*. X* has the same role for FD's as DEP(X) has for MVD's. For each A \in X* we have X \rightarrow A. Thus, each element of X* appears as a singleton set in the dependency basis of X. The dependency basis contains other sets if and only if X* is a proper subset of U.

These remaining sets cover U-X*. We note that since $X \rightarrow X^*$ and $X^* \rightarrow X$, it follows that DEP(X) = DEP(X*) and the dependency basis of X is the same as the dependency basis of X*.

<u>Theorem 1</u>. (Completeness theorem for FDs and MVDs): Let F and M be sets of FDs and MVDs, respectively. For each functional or multivalued dependency that does not belong to $(F,M)^+$ there exists a relation R(U) such that all the dependencies in $(F,M)^+$ are valid in R but the given dependency is not valid in R.

<u>Proof</u>: Let X be the left side of the dependency that is not in $(F,M)^+$. The set X* is a proper subset of U since otherwise every (functional or multivalued) dependency with left side X belongs to $(F,M)^+$. Let W_1, \ldots, W_m $(m \ge 1)$ be the sets in the dependency basis of X that cover U-X*. Thus, X*, W_1, \ldots, W_m form a partition of U. The MVD X \longrightarrow X* $W_1 \ldots W_m$ is in $(F,M)^+$ and, furthermore, if an MVD in $(F,M)^+$ has X as its left side then its right side is a union of a subset of X* and some of the sets W_1, \ldots, W_m .

The relation R(U) is constructed as follows: We choose the set $\{0,1\}$ as the domain of each of the attributes in U. The relation R has 2^{m} rows, one row for each sequence of zeros and ones of length m. In the row corresponding to a sequence $\langle a_{1}, \ldots, a_{m} \rangle$ (where $a_{i} \in \{0,1\}$), each of the attributes in W_i is assigned the value a_{i} ($i = 1, \ldots, m$). Each attribute in X* is assigned the value 1 in all the rows of the relation. For example, if m = 3, then the row corresponding to the sequence $\{0,1,1\}$ has all 1's in the X* columns, all 0's in the W₁ columns, all 1's in the W₂ columns and all 1's in the W₃ columns.

We now want to prove that R satisfies the condition of the theorem. In what follows, we use the inference rules for two different purposes. First, we sometimes show that if some given dependencies are in $(F,M)^+$ then $(F,M)^+$ also contains some other dependency. That we can use the rules for this purpose follows directly from the definition of $(F,M)^+$. Second, we also show that if some dependencies are valid in R then there is another dependency that is valid in R. We can use the rules for this purpose since we have proved that they are valid inference rules for the family of dependencies, i.e., their application to dependencies that are valid in a relation always produces dependencies valid in that relation. We will indicate our intention each time we use the rules.

We now prove the following three claims about the relation R that we have just constructed.

- (1) If the right side of an FD is a nonempty subset of W_i then the FD is valid in R iff its left side intersects W_i.
- (2) Each MVD that has W; as its right side is valid in R.
- (3) If the right side of an MVD is a nonempty proper subset of W_i then the MVD is valid in R iff its left

side intersects W_i.

We first prove one direction of claim 1. For each fixed row of R, all the attributes in W_i have the same value. It follows that every attribute in W_i is functionally dependent in R on every other attribute in W_i (and, by augmentation, on every set that contains such an attribute). Thus we have proved that if the left side of the FD intersects W_i then the FD is valid in R. From this also follows the corresponding direction of claim 3, since every FD is also an MVD.

We now prove claim 2 and the other directions of claims 1 and 3. The relation R is the Cartesian product of its projections $R(W_i)$ and $R(U-W_i)$. It follows immediately that the MVD $\emptyset \longrightarrow W_i$ is valid in R and, by augmentation, $Y \longrightarrow W_i$ is valid in R, for every set Y. This proves claim 2. Now let Y and Z be subsets such that Y is disjoint from W, and Z is a nonempty subset of W,. It follows from the above factorization of R that for each Y-value (Def: For a set of attributes X, an X-value is an assignment of values to the attributes of X from their domains) y, the set $Z_r(y)$ contains two Z-values - a 0assignment and a 1-assignment to the attributes of Z. In particular, the FD Y \rightarrow Z is not valid in R; this concludes the proof of claim 1. If Z is a proper subset of W_i let A be an attribute in W_i - Z. Then $Z_r(ya)$, where ya is a YA-value, contains only a single Z-value since the

attributes of Z must be assigned the same value as the attribute A. Therefore $Z_R(y) \neq Z_R(ya)$ and it follows that the MVD Y \longrightarrow Z is not valid in R. This concludes the proof of claim 3.

We now show that R satisfies the condition of the theorem. First, let f be an FD in $(F,M)^+$. We show that f is valid in R. By Axiom FD6 we can assume that f is of the form $Y \rightarrow B$ where B is a single attribute. Now, if B is in X* then f is clearly valid in R since in R every attribute of X* assumes a single value and is, therefore, functionally dependent on any other attribute. If B is not in X* then it belongs to some W_i . If Y is disjoint from this W_i then from the MVD X $\rightarrow W_i$ and the FD Y \rightarrow B which are both in $(F,M)^+$ it would follow by rule FD-MVD2 that X \rightarrow B is in $(F,M)^+$. This is impossible since B is not in X*. Therefore Y must intersect W_i and Y \rightarrow B is valid in R by claim 1.

Now let g: $Y \longrightarrow Z$ be an MVD in $(F,M)^+$. We show that g is valid in R. We note that $Y \longrightarrow Z \cap X^*$ is valid in R. We will show that, for each i, $Y \longrightarrow Z \cap W_i$ is also valid in R. First, suppose that, for some i, the set $Z \cap W_i$ is either empty or all of W_i . By claim 2 above, $Y \longrightarrow W_i$ is valid in R; as we know, $Y \longrightarrow \emptyset$ is always valid. Next, suppose that, for some i, $Z \cap W_i$ is a nonempty proper subset of W_i . If Y does not intersect W_i we can use augmentation on $Y \longrightarrow Z$ to obtain that $U-W_i \longrightarrow Z$ is in $(F,M)^+$. Since $X \to W_i$ is also in $(F,M)^+$ it follows by MVD3 that $X \to Z - (U - W_i)$, that is $X \to Z \cap W_i$ is in $(F,M)^+$. This is a contradiction to the assumption that W_i is a member of the dependency basis of X. Thus Y must intersect W_i and, by claim 3, $Y \to Z \cap W_i$ is valid in R. We have now shown that, for each i, $Y \to Z \cap W_i$ is valid in R and also so is $Y \to Z \cap X^*$. By taking the union (MVD4) it follows that $Y \to Z$ is valid in R.

Finally, let us consider the dependency (with left side X) which is known not to be in $(F,M)^+$. If it is an FD X \rightarrow Y then Y is not a subset of X*, so Y intersects W_i for some i. By FD6 if X \rightarrow Y is valid in R so is X \rightarrow Y \cap W_i and this contradicts claim 1. (Recall that X* is disjoint from each of the W_i .) Therefore X \rightarrow Y is not valid in R. If the dependency is an MVD X \rightarrow Y, then for some i, Y \cap W_i must be a nonempty proper subset of W_i (else since C \rightarrow Y \cap X* and X \rightarrow Y \cap W_i for each i are in (F,M)⁺ the MVD X \rightarrow Y would be in (F,M)⁺). Now, for this i, X \rightarrow Y \cap W_i is not valid in R by claim 3. Since X \rightarrow W_i is valid in R, if X \rightarrow Y was also valid in R we could apply MVD6 to obtain a contradiction. Thus X \rightarrow Y is not valid in R. # <u>Proposition</u>. The Union rule (i.e., MVD5) cannot be derived from rules MVD1-MVD3.

<u>Proof</u>: Let U = {A,B,C}, G = {A $\rightarrow \rightarrow$ B,A $\rightarrow \rightarrow$ C}, and g = A $\rightarrow \rightarrow$ BC. Clearly g follows from G via MVD5. Let S = {MVD1,MVD2,MVD3}; we claim that $G^S \subseteq G'$, where G^S is the closure of G under S, and G' = {X $\rightarrow \rightarrow$ Y/Y \subseteq X or A \in X and Y-X = B or C}. To prove this, it suffices to show that G' is closed under S, since clearly G is included in G'. We shall consider each rule of S in turn.

- MVD1. Obviously G' contains all dependencies obtainable by reflexivity.
- MVD2. Suppose $X \longrightarrow Y \in G'$ and $X \subseteq W$. Consider two cases.
 - (a) $Y \subseteq X$; then $YZ \subseteq XW$ implies that $XW \longrightarrow YZ \in G'$.
 - (b) $Y \subseteq X$, $A \in X$, and Y-X = B or C. Suppose Y-X = B. Then YZ-XW=(Y-X-W)(Z-X-W) = B-W. If $B \in W$, then $YZ-XW=\emptyset$, which implies that YZ is contained in XW, so $XW \longrightarrow YZ \in G'$. If $B \notin W$, then YZ-XW = Band $A \in XW$ implies that $XW \longrightarrow YZ \in G'$. A similar argument holds for Y-X = C.

MVD3. Suppose $X \rightarrow Y$, $Y \rightarrow Z \in G'$. There are

[†]This proof follows quite closely that of Mendelzon [24].

four cases to consider.

(a) $Y \subseteq X$, $Z \subseteq Y$; then X contains 2-Y, so X \longrightarrow Z-Y ε G'.

(b) $Y \subseteq X$, Y does not contain Z, and A ϵ Y, Z-Y = B or C.

Assume Z-Y = B. If $B \in X$, then Z-Y-X=Ø implies that X contains Z-Y, hence $X \longrightarrow Z-Y \in G'$. If $B \notin X$, then Z-Y-X=B and $A \in X$ (since $A \in Y$ and $Y \subseteq X$); hence $X \longrightarrow Z-Y \in G'$. The argument is similar for Z-Y = C.

(c) $Y \not\subseteq X$, $A \in X$, Y-X=B or C; $Z \not\subseteq Y$, $A \in Y$, Z-Y = B or C.

Suppose Z-Y = B. If $B \in X$, Z-Y-X = \emptyset , hence X \longrightarrow Z-Y as above. If $B \notin X$, Z-Y-X=B and we assumed $A \in X$, so X \longrightarrow Z-Y \in G'. The argument for Z-Y is similar. (d) Y $\notin X$, $A \in X$, Y-X=B or C, and Z \subseteq Y. Now $\overline{Z}-Y = \emptyset$, so X \longrightarrow Z-Y \in G'.

It follows that $A \longrightarrow BC$ is not in G^S , and therefore $A \longrightarrow BC$ cannot be derived from $\{A \longrightarrow B, A \longrightarrow C\}$ using only MVD1-MVD3. # A Membership Algorithm[†]

A membership algorithm which can be applied well to real world applications is presented below. In this procedure EQN is a set of numbers associated with FDs in F. Membership Algorithm

An Algorithm to decide if an FD is in the closure of a set of FDs INPUT : A set of FDs F in canonical form; and an FD g, whose left and right sides are LM and RM, where |RM| = 1."YES", if $q_{\epsilon}F^{\dagger}$, and "NO" otherwise. OUTPUT: $X = \emptyset;$ procedure MEMBER (LM, RM, EQN) do for all $i \in EQN$; $\underline{if} RIGHT_{\underline{f}} = RM$ <u>then</u> $\underline{if} \ \underline{LEFT}_{f_i} \subseteq \underline{LM} \ \underline{then} \ \underline{return}('YES');$ else [<u>do</u> for all d ε (LEFT_f, - LM); then go to end d-loop; if dεX SUB \leftarrow MEMBER(LM, d, EQN - {i}); if SUB = 'NO' then $[EQN = EQN - {i};$ exit;

[†]This work was done by Dr. John C. Thompson and the author, and was presented at the Sixth Annual Computer Science Conference at Detroit, Michigan in February, 1978.

```
end d-loop;
if SUB = 'YES' then [X = X [Jd; return('YES')];
]
end i-loop;
return('NO');
end MEMBER;
```

Proof of termination

Theorem 2 The membership algorithm halts.

<u>Proof</u>: The outer loop is finite because EQN is a finite set. At each iteration of the inner loop, either the value of SUB is 'YES' for all d's (including those created by recursion) in which case the algorithm terminates, or at least one FD is crossed off the set EQN. Now, since the number of elements in EQN is finite, the algorithm ultimately halts. #

Proof of correctness

Theorem 3 Upon termination of the membership algorithm, the output is "YES" if and only if $g \in F^+$.

<u>Proof</u>: At the first call of the procedure, the right sides of all FDs in F are checked until an FD f_i whose right side is the same as the right side of g is found (if no such an FD is found, g cannot be in the closure of F, and the algorithm correctly returns "NO" as its output). If the left side of f, is a subset of the left side of g, then clearly g is in the closure of F, and this is performed by the third statement of the procedure. If the left side of the f; is not a subset of the left side of q, then the procedure is recursively called for each element d in the set difference of these two left sides (i.e., LEFT_{f_s} - LM). Therefore, the subsequent procedure call is similar to the first call, and it determines if there is any FD $\epsilon(F - \{f_i\})$ with the right side of d. . If found, again its left side is checked and the same procedure is continued until we get to a point in which all d's are replaced by the left sides of their corresponding dependencies (or the subsequent left sides), and then this set is either a subset of the left side of g, or it is not. The former case indicates that $g \in F^+$, and the inner loop of the algorithm does this job correctly. On the other hand, the latter condition states that f; cannot be used to derive g from the set F, and thus the algorithm checks for the other FDs in F. Finally, if all FDs of F are checked (when the outer loop is exhausted) with the same situation as mentioned above, the algorithm returns "NO" and terminates. #

Speed analysis

The time analysis of the algorithm can be better explained by considering the worst situation that may

occur. Suppose there are n FDs as follows:

$$f_{1}: A_{1}, A_{2}, \dots, A_{m1} \rightarrow B_{1}$$

$$f_{2}: A_{1}, A_{2}, \dots, A_{m2} \rightarrow B_{2}$$

$$\vdots$$

$$f_{n}: A_{1}, A_{2}, \dots, A_{mn} \rightarrow B_{n}$$

and suppose we want to know if an FD $\alpha \rightarrow \beta$ is in the closure of the above set. The algorithm first checks for β in the left side of FDs. Suppose there is an f, in the set such that $B_i = \beta$ (if no such an FD exists, the algorithm terminates in O(n)). Now, the worst situation happens if LEFT_f, $\Omega \alpha = \emptyset$, and the algorithm has to search for all $A_k \in LEFT_{f_i}$, that is, for m_i attributes in (n-1)dependencies. Notice that, even if $m_i > (n-1)$, the inner loop may not be executed for more than (n-1) times. This is because the nth element of the set LEFT f. cannot be found in the right side of the given FDs (assuming the first (n-1) elements have already been found), causing an exit from the loop. Now the worst situation occurs if the first element can be found after searching the whole set, i.e., (n-1) FDs. Since this dependency is removed from the set, and the element found is added to the set X (i.e., a set which keeps track of these attributes to avoid re-searching for them in the subsequent procedure calls), then for the next phase (i.e., searching for the attributes on the left side of the dependency just

removed) we have (n-2) dependencies to search. Continuing a similar procedure we get a total of (n-1)+(n-2)+... + 1, and hence the overall time for the algorithm is below $0(n^2)$.

As the reader will notice, the algorithm first checks the right side of the FDs. And since all FDs are in canonical form (i.e., they have only one attribute on their left sides), this test can be efficiently implemented. The algorithm also checks the left side of the dependencies, but this test has been minimized. In other words, the left side of a dependency is checked only if its right side is appropriate (in the sense that this dependency might be used for derivation of the FD of interest).

In addition, in this algorithm, only those FDs that might be used for the derivation purpose will be tested, and the other dependencies will not be checked at all.