

Extending the E-R Concepts to Capture Natural Language Semantics for Database Access

Frank Shou-Cheng Tseng

Department of Computer Science and
Information Engineering
National Chiao Tung University
Hsinchu, Taiwan
dcp77806@csunix.csie.nctu.edu.tw

Arbee L. P. Chen

Department of Computer Science
National Tsing Hua University
Hsinchu, Taiwan
alpchen@cs.nthu.edu.tw
Fax: 886-35-723694

Wei-Pang Yang

Department of Computer Science and
Information Engineering
National Chiao Tung University
Hsinchu, Taiwan
wpyang@twncu01.bitnet

Abstract

Research on accessing databases in natural language usually employs an intermediate form for the mapping process from natural language to database languages. However, much effort is needed to bridge the gap between the existing intermediate forms and the database languages. In this paper, we extend the E-R concepts to capture the natural language semantics and develop a logical form to represent the natural language queries. The logical form can be efficiently transformed into relational algebra for query execution. The whole process provides a clear and natural framework for processing natural language queries to retrieve data from database systems.

1 Introduction

In database management systems, data are retrieved via a well-defined query language. Although some query languages, say SQL (Structural Query Language), can be very powerful, users may suffer greatly from their complex usage. Research in natural language processing for DBMS interface attempts to ease this complexity by freeing users from knowing the exact database structure and learning the query language.

An intermediate form is usually used for mapping natural language constructs onto the underlying database schema. This intermediate form can be a universal relation [16] as used in FRED [10][14]. It can also be a hierarchy structure constructed from the information about the database. TEAM [9] and QPROC [17] are two example systems of this type. Besides, the TQA [6] system use a semantic net model as its intermediate representation.

However, Ullman has pointed out that interpreting queries over a universal relation is a difficult task [16]. Other intermediate forms suffer from the bias to natural language constructs and much effort is needed to translate the intermediate forms into the database languages.

In this paper, we study the inter-relationship between the natural language constructs of a query and the E-R conceptual schema [3]. Chen [4] has pointed out that the basic constructs of English sentences can be mapped onto E-R schemas in a natural way. Chen [4] has studied 11 rules for translating information requirements, which are originally documented in English, into database schemas in terms of E-R diagrams. In comparison, we propose an approach to map the natural language queries into relational algebra through the E-R schema.

1.1 Review of E-R Model

The Entity-Relationship model [3] adopts the view that the real world consists of *entities* and *relationships* among entities. An entity is a 'thing' which can be distinctly identified. A relationship is an association among entities. The E-R model uses the concepts of *entity set*, *relationship set* and *value set*. More details about these concepts can be found in [3].

The structure of a database organized according to the E-R data model can be depicted by an Entity-Relationship Diagram (ERD). In an ERD, an entity set is represented by a rectangular, and a relationship set is represented by a diamond, labeled by their associated names. The entity sets that participate in a relationship set are indicated by edges with their corresponding semantic roles attached. Moreover, a value set is represented as an oval labeled with the attribute defined on it.

Suppliers			
sno	sname	status	city
S1	Smith	20	London
S2	Jones	10	Paris
S3	Blake	30	Paris
S4	Clark	20	London
S5	Adames	30	Taipei

Shipments		
sno	pno	qty
S1	P1	300
S1	P2	200
S1	P3	400
S1	P4	200
S1	P5	100
S1	P6	100
S2	P1	300
S2	P2	400
S3	P2	200
S4	P2	200
S4	P4	300
S4	P5	400

Parts				
pno	pname	color	weight	city
P1	Nut	Red	12	London
P2	Bolt	Green	17	Paris
P3	Screw	Blue	17	Rome
P4	Screw	Red	14	London
P5	Cam	Blue	12	Paris
P6	Cog	Red	19	London

Figure 1: The Suppliers-and-Parts Database.

Information about an entity set can be organized into a relation which is named *entity relation*. Similarly, a *relationship relation* can be constructed by collecting the attributes of the relationship set with the primary keys of the associated entity relations. For example, in Figure 1, the relations Suppliers and Parts of the Suppliers-and-Parts database described in [8] are *entity relations*, while Shipments is a *relationship relation*. Our examples will be based on this database hereafter.

1.2 Overview of Natural Language Processing

Language understanding process is commonly divided into three stages. First, the sentence is parsed according to the predefined grammar, then the semantic roles are built and finally these semantic roles are mapped onto the specific objects in the real world.

Augmented-Transition-Net (ATN) [19] is traditionally used to parse a natural language sentence. By following a method for expressing grammars in logic due to Kowalski [11], Pereira and Warren [13] have developed a clearer and more powerful formalism named Definite Clause Grammar (DCG). McCord [12] contributes to the syntactic analysis and semantic interpretation of natural languages in the framework of logic programming.

Winston [18] described a variety of constraints to help establish semantic roles in a sentence. These semantic roles reveal how the nouns are related to the verb. Consider the sentence "Andy moves the dirty table into the messy office." It is parsed into a noun phrase and a verb phrase. The verb phrase consists of a noun phrase and a prepositional phrase which consists of another noun phrase.

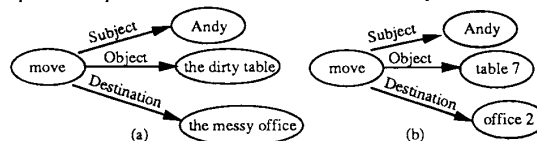


Figure 2 (a): The Semantic Roles for "Andy moves the dirty table into the messy office."
(b): Mapping Semantic Roles onto Specific Objects.

The semantic roles of the sentence are shown in Figure 2(a). It

indicates that the *verb* is move, the *subject* is Andy, the *object* is the dirty table, and the *destination* is the messy office. Finally, these semantic roles can be mapped onto specific objects in the real world as Figure 2(b) illustrates.

Our approach follows these three stages and we focus on the mapping from semantic roles to an E-R schema. We develop a logical form to represent the mapping result and it can be efficiently transformed into the relational algebra for query execution.

The remainder of this paper is organized as follows. Section 2 describes the processing model and the mapping process. Section 3 devotes to the logical form constructs which can be used to represent the mapping result of a user query. Query transformation process which transforms logical forms into the relational algebra is described in Section 4. Finally, we conclude and suggest future research in Section 5.

2 The Processing Model and the Mapping Process

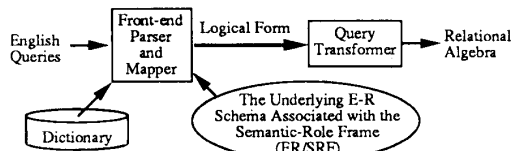


Figure 3: The Processing Model.

2.1 The Processing Model

The processing model is shown in Figure 3. There is a front-end parser/mapper to parse the English queries and map them onto the underlying E-R schema. The parsing and mapping process may refer to

1. the dictionary and
2. the underlying E-R schema associated with the semantic-role frame (ER/SRF).

After the parsing phase, the query is decomposed into

1. *semantic roles*, each of which is composed of a *headnoun* and some *modifiers* (a headnoun is the main noun in a noun phrase; for example, "the London suppliers" has a headnoun 'suppliers' and the other noun 'London' is a modifier which modifies the headnoun) and
2. the *verb* that relates these semantic roles.

Each of the semantic roles is mapped onto an *entity relation* and its headnoun and modifiers are mapped onto the corresponding attributes of that entity relation based on ER/SRF (to be discussed in Section 2.1.2). The verb that relates these semantic roles is mapped onto the relationship relation that associates these entity relations. Figure 4 illustrates this mapping mechanism.

We develop a logical form in Section 3 to represent the mapping result. After the logical form is generated, it is passed to the query transformer to produce the query in relational algebra.

Our processing model makes the following assumptions.

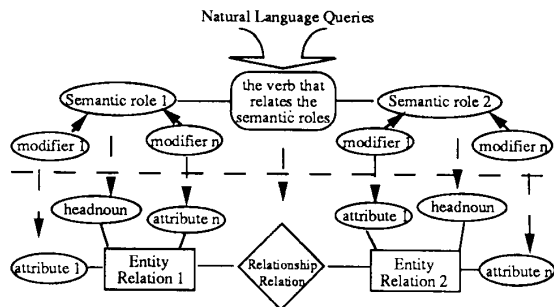


Figure 4: The Schema Mapping Mechanism.

1. In addition to the original attributes, each entity relation is augmented with a *surrogate*. Any reference to this surrogate is interpreted as a reference to the corresponding relation. For example, in the Suppliers-and-Parts database, we add the extra attributes *sno* and *pno* as the surrogates of Suppliers and Parts, respectively.

2. We also augment each relationship relation with a surrogate which is composed of the surrogates of the involved entity relations. For example, the surrogate of the relationship relation Shipments is formed by grouping *sno* and *pno* into (*sno*, *pno*).

2.1.1 The Dictionary

The linguistic knowledge that enables the semantic roles of a verb to be mapped onto the correct attributes of the associated relationship relation is stored in the dictionary. Notice that some of these attributes are surrogates of the entity relations involved in the relationship. For the Suppliers-and-Parts database, the knowledge for the verb 'supply' would be like

Relation	Semantic Roles	the Corresponding Attribute
Shipments	Subject	sno (the surrogate of Suppliers)
	Object	pno (the surrogate of Parts)
	with_object	qty

Other information such as the synonyms of the entity sets with the corresponding relations and all the domain values with the corresponding attributes are also needed to be stored in the dictionary.

2.1.2 The Schema and the Semantic-Role Frame Represented in ERD

For each entity relation, attributes are identified as the *headnoun* and the *modifiers* that modify the headnoun. For example, the attributes corresponding to the headnouns of Suppliers and Parts are *sname* and *pname*, respectively. Other attributes are identified to be the corresponding modifiers. The above information will be encoded into the original ERD. For example, the ERD and the semantic-role frame representing the Suppliers-and-Parts database are shown in Figure 5.

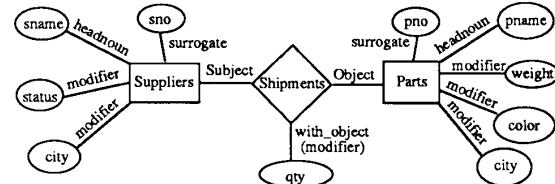


Figure 5: The ERD and the Semantic-Role Frame of the Suppliers-and-Parts Database.

2.2 Description of the Mapping Process

The mapping process can be considered as the counterparts presented in [4].

2.2.1 Mapping Verbs onto Relationship Relations

We can map the verb of a query sentence onto the corresponding relationship relation and its semantic roles onto the corresponding attributes. The query "Does Smith supply nuts with quantity 300 ?", for instance, has the semantic roles which can be mapped onto the attributes as depicted in Figure 6. The attributes *sno* and *pno* refer to the entity relations Suppliers and Parts, respectively. That tells us 'Smith' and 'nuts' will be associated with the headnouns of the entity relations Suppliers and Parts, respectively. Moreover, '300' will be mapped onto the attribute *qty* of the relationship relation Shipments.

Sentence	Does	Smith	supply	nuts	with quantity 300 ?
Semantic roles		Subject	verb	Object	with_object
		sno (attribute) (surrogate)		pno (attribute) (surrogate)	qty (attribute)
		Shipments (relationship relation)			

Figure 6: The Mapping for "Does Smith supply nuts with quantity 300 ?"

Note that the verb 'be' is treated differently, for it is a linking verb, which is followed by a subject complement which describes the subject. Linking verbs do not transfer action; rather, they join the subject and the complement. For example,

Smith is a London supplier.

in which *supplier* describes *Smith*. Because the verb 'be' does not transfer action, we need not associate it with a *relationship relation*. The above sentence is mapped onto the entity relation Suppliers (and the attributes *sname* and *city*, see Section 2.2.2).

Note also that an imperative mood sentence is always used for issuing command. The leading verb does not transfer action. For example,

List the suppliers located in London.

is mapped onto the entity relation Suppliers (and attribute *city*).

2.2.2 Mapping Noun Phrases onto Entity Relations and Its Modifiers onto Attributes

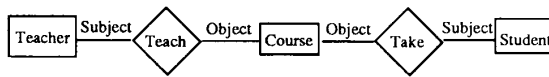
As we have shown in the previous section, the semantic roles of a verb in a sentence can be mapped onto those surrogates of involved entity relations. Those entity relations actually interpret noun phrases of the sentence.

However, the noun phrases may consist of more than one noun modifiers (which can be adjectives or nouns) or relative clauses that modify the headnoun. For example, "the London suppliers" and "the red parts" are two noun phrases and have the noun modifiers 'London' and 'red', respectively. These noun modifiers can be interpreted into the attributes of the corresponding entity relation. According to the dictionary, the interpretation of these noun modifiers described above can be mapped as Figure 7 illustrates. Notice that some relationships

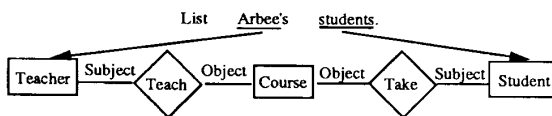
Noun Phrase	Corresponding Relation	Interpretation of the Noun Modifier
the London Suppliers	Suppliers	city = "London"
the red parts	Parts	color = "red"

Figure 7: Mapping Noun Modifiers onto Attributes.

can be queried in possessive form like "List A's B." (or "List the B of A.") In this case, there is no verb to be mapped onto the relationship relation. But, by referring to the schema and the semantic-role frame represented in ERD, we may first map the noun phrases (A and B) onto the corresponding entity relations according to the ERD. This path is obtained as the mapping result. For example, if we have the following schema



then the query "List Arbee's student." can be mapped as follows.



And the relations on the path (Teacher, <Teach>, Course, <Take>, Student) are our mapping result, where relationship relations are enclosed by angle brackets. Note that this query is equivalent to "List the students who take the courses taught by Arbee."

3 The Logical Form

In this section, we develop a logical form for representing the results of mapping from English queries onto an E-R schema. It can be easily translated into the relational algebra for query execution.

3.1 The Extension of E-R Diagram for the Logical Form

Based on the constructs of E-R diagram, we develop a logical form by extending the constructs of E-R diagram. The logical form takes into account the following conditions of an English query.

1. The representation for the negative and positive forms of a verb,

2. The representation for modifiers,
3. The representation for conjunctives 'and' and 'or', and
4. The representation for the word 'all'.

These conditions are explained in the following subsections.

3.1.1 The Representation for the Negative and Positive Forms of a Verb

A verb is usually mapped onto a relationship relation, which is depicted as a diamond in the E-R diagram. But the verb in a query may be issued in negative form. For example, a user may issue a query in negative form like "List the suppliers who *do not supply* nuts." We extend the diamond representation of E-R diagram to represent both the negative and positive form of a query as Figure 8 shows. In posi-

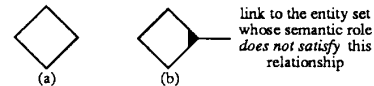


Figure 8: (a) The Relationship Representation in Positive Form (b) The Relationship Representation in Negative Form.

tive case, we represent the relationship as before; in negative case, we represent the relationship as a diamond except that there is a black triangle that links to the entity set whose semantic role *does not satisfy* the relationship. In the previous query, the black triangle links to the entity set Suppliers.

3.1.2 The Representation for Modifiers

After the headnoun and modifiers of an entity set are recognized, we respectively associate them with the corresponding attributes and form them into predicates (attribute θ constant, $\theta \in \{>, <, =, \neq, \geq, \leq\}$). These predicates are represented by oval nodes in our logical form. For example, "List the suppliers who supply red parts." has the logical form shown in Figure 9. We define the predicate "attribute = ?" as a *pseudo predicate*; it represents the target attribute which is to be output to the user.

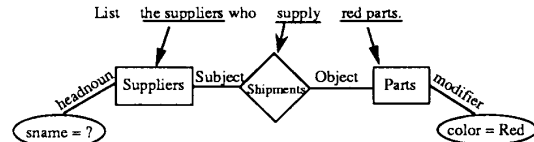


Figure 9: The Logical Form of "List the suppliers who supply red parts."

3.1.3 The Representation for Conjunctives 'And' and 'Or'

The conjunctive 'and' is sometimes interpreted to be a disjunction (logical 'or'). For example, "List the red parts and the blue parts." is equivalent to "List the red parts or the blue parts." But we will not explore such phenomenon in this paper. We assume ambiguities are solved in parsing phase.

From the point of view of E-R model, 'and' and 'or' can be used to conjunct the following three cases.

1. *Modifier and Modifier*. In this case, the oval nodes that correspond to these modifiers are linked pairwise by edges labeled 'A' and 'V' for conjunctives 'and' and 'or', respectively. For example, "List the parts with *color blue* and (or) *located in London*." has the logical form shown in Figure 10(a).
2. *Entity Set and Entity Set*. In this case, a conjunctive conjuncts two noun phrases. The rectangle nodes that correspond to the nouns are linked pairwise by edges labeled 'A' and 'V' for conjunctives 'and' and 'or', respectively. For instance, "List the suppliers who *supply red parts* and (or) *nuts*." has the logical form shown in Figure 10(b).
3. *Relationship and Relationship*. In this case, a conjunctive conjuncts two verb phrases. The diamond nodes that correspond to the verbs are linked pairwise by edges labeled 'A' and 'V' for conjunctives 'and' and 'or', respectively. For instance, "List the suppliers who *supply red parts* and (or) *supply nuts*." has the logical form shown in Figure 10(c). Note that the following example does not fall into this case.

List the suppliers who supply nuts and are located in London.

Although the 'and' conjuncts two verbs, the verb "are located in" is used to modify the suppliers instead of performing action. This sentence is equivalent to "List the London suppliers who supply nuts." and there is no conjunction at all.

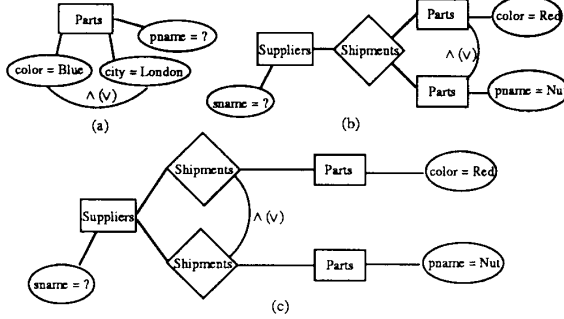


Figure 10: (a) The Logical Form for "List the parts with color blue and (or) located in London." (b) The Logical Form for "List the suppliers who supply red parts and (or) nuts." (c) The Logical Form for "List the suppliers who supply red parts and (or) supply nuts."

3.1.4 The Representation for the Word 'All'

Consider the following examples,

- List the suppliers who supply *all* red parts.
- List the suppliers who supply red parts.

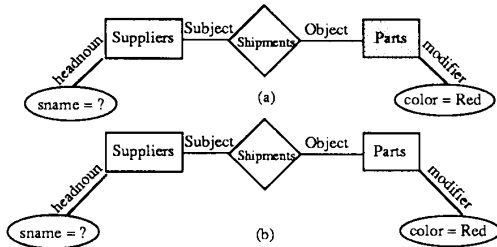


Figure 11: (a) The Logical Form for "List the suppliers who supply all red parts." (b) The Logical Form for "List the suppliers who supply red parts."

The answer of the first example is the suppliers who supply *all* the red parts, but the second example is the suppliers who supply *any* red parts. Therefore, the answer of (1) is always contained in that of (2). We use a shadow rectangle to represent the entity set whose semantic role is preceded by 'all'. Otherwise, a blank rectangle is used (see Figure 11). Note that the semantic meaning of "List *all* the suppliers who supply *all* red parts." is equivalent to that of (1). The first 'all' has no effect on this query.

3.2 Definition of the Logical Form

From the previous discussion, a logical form for a query Q can be denoted by $LF(Q) = (N, E, f_N, f_E)$, where

- N is a set of nodes which can be further classified into the sets N_r , N_e , and N_p (i.e. $N = N_r \cup N_e \cup N_p$), where
 - N_r is the set of diamond nodes representing *relationship relations*,
 - N_e is the set of rectangle nodes representing *entity relations*,
 - N_p is the set of oval nodes representing *predicates* which is of the form "attribute θ constant", $\theta \in \{>, <, =, \neq, \geq, \leq\}$.
- E is a set of edges, which can be further classified into the sets $E_{(r,e)}$, $E_{(r,p)}$, $E_{(e,p)}$, $E_{(p,p)}$, $E_{(e,e)}$, and $E_{(r,r)}$ (i.e. $E = E_{(r,e)} \cup E_{(r,p)} \cup E_{(e,p)} \cup E_{(p,p)} \cup E_{(e,e)} \cup E_{(r,r)}$), where
 - $E_{(r,e)} \subseteq N_r \times N_e$. An edge $(r, e) \in E_{(r,e)}$ is said to join the diamond node r and the rectangle node e .

- $E_{(r,p)} \subseteq N_r \times N_p$. An edge $(r, p) \in E_{(r,p)}$ is said to join the diamond node r and the oval node p .
- $E_{(e,p)} \subseteq N_e \times N_p$. An edge $(e, p) \in E_{(e,p)}$ is said to join the rectangle node e and the oval node p .
- $E_{(p,p)} \subseteq N_p \times N_p$. An edge $(p_1, p_2) \in E_{(p,p)}$ is said to join the oval nodes p_1 and p_2 .
- $E_{(e,e)} \subseteq N_e \times N_e$. An edge $(e_1, e_2) \in E_{(e,e)}$ is said to join the rectangle nodes e_1 and e_2 .
- $E_{(r,r)} \subseteq N_r \times N_r$. An edge $(r_1, r_2) \in E_{(r,r)}$ is said to join the diamond nodes r_1 and r_2 .

3. f_N is a set of mappings, $f_N = \{f_{N_r}, f_{N_e}, f_{N_p}\}$, where

- $f_{N_r} : N_r \rightarrow R_N$, where R_N is the set of the relationship relation names labeled on $r_i, \forall r_i \in N_r$.
- $f_{N_e} : N_e \rightarrow E_N \times \{\forall, \exists\}$, where E_N is the set of the entity relation names labeled on $e_i, \forall e_i \in N_e$. $\{\forall, \exists\}$ represents the cases 'all' (\forall) and 'any' (\exists) as addressed in Section 3.1.4.
- $f_{N_p} : N_p \rightarrow P$, where P is the set of the predicates labeled on $p_i, \forall p_i \in N_p$.

4. $f_E = \{f_{E_{r,e}}, f_{E_{r,p}}, f_{E_{e,p}}, f_{E_{p,p}}, f_{E_{e,e}}, f_{E_{r,r}}\}$ is a set of mappings, where

- $f_{E_{r,e}} : E_{(r,e)} \rightarrow S_r \times \{P, N\}$, where S_r is the set of the labels of $(r_i, e_j), \forall (r_i, e_j) \in E_{(r,e)}$, which represent the *semantic roles* of the rectangle nodes e_j . $\{P, N\}$ represents the positive (P) and negative (N) cases discussed in Section 3.1.1.
- $f_{E_{r,p}} : E_{(r,p)} \rightarrow VM$, where VM is the set of labels of $(r_i, p_j), \forall (r_i, p_j) \in E_{(r,p)}$, which represent the *verb modifiers* of the verb corresponding to the diamond nodes r_i .
- $f_{E_{e,p}} : E_{(e,p)} \rightarrow \{\text{headnoun}, \text{modifier}\}$.
- $f_{E_{p,p}} : E_{(p,p)} \rightarrow \{\wedge, \vee\}$, where ' \wedge ' and ' \vee ' are the edge labels representing conjunction and disjunction, respectively.
- $f_{E_{e,e}} : E_{(e,e)} \rightarrow \{\wedge, \vee\}$.
- $f_{E_{r,r}} : E_{(r,r)} \rightarrow \{\wedge, \vee\}$.

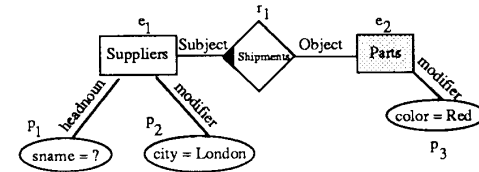


Figure 12: The Logical Form for Example 3.1.

Example 3.1 For the query $Q =$ "List the London suppliers who *do not* supply *all* red parts.", the logical form $LF(Q) = (N, E, f_N, f_E)$ as represented in Figure 12 is formally specified as follows:

- $N = \{r_1, e_1, e_2, p_1, p_2, p_3\} = N_r \cup N_e \cup N_p$, where $N_r = \{r_1\}$, $N_e = \{e_1, e_2\}$, and $N_p = \{p_1, p_2, p_3\}$.
- $E = \{(r_1, e_1), (r_1, e_2), (e_1, p_1), (e_1, p_2), (e_2, p_3)\} = E_{(r,e)} \cup E_{(r,p)} \cup E_{(e,p)} \cup E_{(e,e)} \cup E_{(r,r)}$, where $E_{(r,e)} = \{(r_1, e_1), (r_1, e_2)\}$, $E_{(e,p)} = \{(e_1, p_1), (e_1, p_2), (e_2, p_3)\}$. $E_{(r,p)}$, $E_{(p,p)}$, $E_{(e,e)}$, and $E_{(r,r)}$ are all empty sets.
- $f_N = \{f_{N_r}, f_{N_e}, f_{N_p}\}$, where
$$f_{N_r}(r_1) = (\text{Shipments})$$

$$f_{N_e}(e_1) = (\text{Suppliers}, \exists), f_{N_e}(e_2) = (\text{Parts}, \forall),$$

$$f_{N_p}(p_1) = "sname = ?", f_{N_p}(p_2) = "city = London", \text{ and}$$

$$f_{N_p}(p_3) = "color = red".$$
- $f_E = \{f_{E_{r,e}}, f_{E_{e,p}}\}$, where

$f_{Ere}((r_1, e_1)) = (\text{Subject}, N), f_{Ere}((r_1, e_2)) = (\text{Object}, P),$
 $f_{Esp}((e_1, p_1)) = \text{'headnoun'}, f_{Esp}((e_1, p_2)) = \text{'modifier'},$
 and $f_{Esp}((e_2, p_3)) = \text{'modifier'}$. \square

Note that this definition of the logical form can be extended if more natural language constructs are to be taken into account in the mapping process. For example, we may add the constructs for aggregation functions (i.e. Max, Min, Avg, Sum, and Count).

4 Transforming a Logical Form into Its Relational Algebra

We select relational algebra as our target for the following reasons.

1. For the consideration of query optimization.
2. The relational algebra can be further transformed into other query languages (e.g. SQL or QUEL [7]) either for portability consideration [14][15] or for distributed database retrieval [5].

Recall that a *pseudo predicate* is a predicate of the form "attribute = ?". The *target attributes* of a query Q , denoted $T(Q)$, are defined as the set of the attributes involved in all the pseudo predicates of the logical form $LF(Q)$. The i -th component of an n -tuple $t = (c_1, c_2, \dots, c_n)$ is denoted $\pi_i(t) \equiv c_i$. In Example 3.1, $\pi_1(f_{Ne}(e_1)) = \text{'Suppliers'}$ and $\pi_2(f_{Ne}(e_1)) = \text{'\exists'}$.

In the following, the notations used in [16] will be adopted, in which $\sigma, \pi, \bowtie, \cup, \cap, -$, and \bowtie represent selection, projection, division, join, union, intersection, difference, and semijoin, respectively. The transformation process is discussed based on whether the logical form contains diamond node or not. In Section 4.1, we discuss the case where the logical form contains no diamond node. In Section 4.2, we devote to the cases where the logical form contains one or more diamond nodes.

4.1 The Transformation Process for a Logical Form with No Diamond Node

The query transformation process for a logical form $LF(Q)$ is

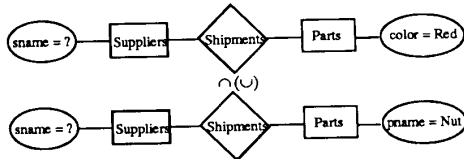
1. For the entity relation, use the predicates, except for pseudo predicates, to restrict it.
2. Project $T(Q)$.

That is, $LF(Q)$ will be transformed into $\pi_{T(Q)}(\sigma_{P_i}(\pi_1(f_{Ne}(e))))$, where e is the single rectangle node, and P_i is the compound predicate of the relation $\pi_1(f_{Ne}(e))$.

4.2 The Transformation Process for a Logical Form Containing One or More Diamond Nodes

4.2.1 The Transformation Process for a Logical Form Containing One Diamond Node

Without loss of generality, we assume that the logical form has no edge of $E(e, e)$. If there are edges of $E(e, e)$ then the logical form can be decomposed according to these edges and the answer is the union/intersection of the results of the sub-logical forms. For example, the logical form in Figure 10(b) can be decomposed into the following two sub-logical forms.



A relational algebra expression represents an execution order of the query. The transformation of a query Q to a relational algebra expression is essentially to determine an algebraic order from its logical form.

The relationship set relates a set of entity sets. That is, the surrogate values of the entity relations can be used to semi-join (i.e. to restrict [1]) the surrogate values of the relationship relation to produce the *answer surrogate values* (which will be used to compute the answer). Before performing the semi-joins, different conditions such as the one-to-all relationship (Section 3.1.4), the negative form relationship (Section 3.1.1), and the predicates (Section 3.1.2) can be evaluated to reduce the surrogate values of the relationship set and entity sets.

The predicates can be directly applied to restrict the corresponding relationship relation or entity relations first. A one-to-all relationship can be implemented by a division operation [8]. Moreover, a negative form relationship can be implemented by computing the positive form relationship followed by a set difference operation.

The transformation process can therefore be stated in the following five phases.

1. For each entity and relationship relations, use their respective predicates, except for pseudo predicates, to restrict the corresponding entity and relationship relations. That is, the transformation first produces $\sigma_{P_i}(\pi_1(f_{Ne}(e_i)))$ and $\sigma_{P_j}(\pi_1(f_{Nr}(r_j)))$, $\forall e_i \in N_e$ and $\forall r_j \in N_r$, where P_i and P_j are the compound predicates of the relations $\pi_1(f_{Ne}(e_i))$ and $\pi_1(f_{Nr}(r_j))$, respectively.

2. Project the surrogates of all relations. Let S_i denote the set of surrogate values of the restricted entity relation RE_i obtained in Phase (1), and (S_1, S_2, \dots, S_n) denotes the set of surrogate values of the restricted relationship relation $RR_{(1,2,\dots,n)}$ which associates with the entity relations $RE_1, RE_2, \dots,$ and RE_n .

3. For all entity relations which correspond to shadow rectangle nodes (i.e., there exist one-to-all relationships), say S_1, S_2, \dots, S_k , define $\langle S_{k+1}, S_{k+2}, \dots, S_n \rangle \equiv$

$$(\dots(((S_1, S_2, \dots, S_n) \div S_1) \div S_2) \dots) \div S_k$$

Also, define $[S_1, S_2, \dots, S_n] \equiv$

$$(S_1, S_2, \dots, S_n) \underset{\bowtie}{S_{k+1}, S_{k+2}, \dots, S_n} \langle S_{k+1}, S_{k+2}, \dots, S_n \rangle.$$

4. If the diamond node corresponding to the relationship relation has a black triangle linked to an entity relation (i.e., there exist negative form relationships), say S_i , then the set of the answer surrogate values can be defined as $\subset S_1, S_2, \dots, S_n \supset \equiv$

$$(S_1, S_2, \dots, S_n) \underset{\bowtie}{S_i} (S_i - \pi_{S_i}([S_1, S_2, \dots, S_n])),$$

else

$$\subset S_1, S_2, \dots, S_n \supset \equiv [S_1, S_2, \dots, S_n].$$

5. Since $\subset S_1, S_2, \dots, S_n \supset$ contains the answer surrogate values, we can join these values with the surrogate values of all relations and project on the target attributes to get the answer. That is, the answer of the query can be produce by

$$\pi_{T(Q)}(\subset S_1, S_2, \dots, S_n \supset \underset{\bowtie}{S_1, S_2, \dots, S_n} RR_{(1,2,\dots,n)})$$

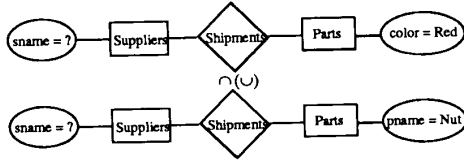
$$\underset{\bowtie}{S_1} RE_1 \underset{\bowtie}{S_2} RE_2 \dots \underset{\bowtie}{S_n} RE_n.$$

However, we can often eliminate unnecessary join operations. If, for example, the attributes of RE_i is not contained in $T(Q)$ then the join on RE_i can be eliminated.

4.2.2 The Transformation Process for a Logical Form Containing More Than One Diamond Node

If a logical form contains more than one diamond node linked by edges of $E(r, r)$ then, without loss of generality, it can be decomposed into sub-logical forms according to the edges. The sub-logical forms can be separately transformed by the process presented in Section 4.2.1.

The final answer is the union/intersection of the results of these sub-logical forms. For example, the logical form in Figure 10(c) can be decomposed into the following two sub-logical forms.



The result of this decomposition is equivalent to that of Figure 10(b). This is because the query corresponding to Figure 10(b) has the same semantic meaning as the query corresponding to Figure 10(c).

Note that if the logical form containing more than one diamond node but these nodes are not linked by edges of $E_{(r,r)}$, then we can transform it as follows.

1. Restrict all the relations by the respective predicates, except for pseudo predicates,
2. Project the surrogates of the restricted relationship relations,
3. Join these surrogates to get the answer surrogate values,
4. Perform necessary joins and project those attributes required by the user.

Such a case occurs when the logical form is constructed from a query which represents the relationship in possessive form (refer to Section 2.2.2).

In the following, we follow Example 3.1 (Figure 12) to show the transformation process. Refer to Figure 1 for this example.

Example 4.1 In Figure 12, we depict the logical form $LF(Q) = (N, E, f_N, f_E)$, where $Q =$ "List the London suppliers who do not supply all red parts.", and we have presented $LF(Q)$ in Example 3.1 in detail. The transformation process is now explained as follows. Notice that we also evaluate the algebraic operations in the process.

1. After restricting both entity relations Suppliers and Parts, we obtain

$$RE_1 \equiv \sigma_{city=London}(Suppliers) \text{ and } RE_2 \equiv \sigma_{color=red}(Parts).$$

Note that there is no restriction on Shipments. Therefore, $RR_{(1,2)} \equiv Shipments$.

2. Perform the projections on surrogates, we get the surrogates as follows.

$$S_1 \equiv \pi_{sno}(RE_1) = \pi_{sno}(\sigma_{city=London}(Suppliers)) = \{S1, S4\}, \\ S_2 \equiv \pi_{pno}(RE_2) = \pi_{pno}(\sigma_{color=red}(Parts)) = \{P1, P4, P6\}, \\ \text{and } (S_1, S_2) \equiv \pi_{sno,pno}(RR_{(1,2)}) = \pi_{sno,pno}(Shipments) \\ = \{(S1, P1), (S1, P2), (S1, P3), (S1, P4), (S1, P5), (S1, P6), \\ (S2, P1), (S2, P2), (S3, P2), (S4, P2), (S4, P4), (S4, P5)\}.$$

3. $\langle S_1 \rangle \equiv (S_1, S_2) \div S_2$ and

$$[S_1, S_2] \equiv (S_1, S_2) \stackrel{S_1}{\bowtie} \langle S_1 \rangle = (S_1, S_2) \stackrel{S_1}{\bowtie} ((S_1, S_2) \div S_2) \\ = (S_1, S_2) \stackrel{S_1}{\bowtie} (\{S1\}) \\ = \{(S1, P1), (S1, P2), (S1, P3), \\ (S1, P4), (S1, P5), (S1, P6)\}.$$

4. $\subset S_1, S_2 \supset \equiv (S_1, S_2) \stackrel{S_1}{\bowtie} (S_1 - \pi_{s_1}(\{S_1, S_2\})) \\ = (S_1, S_2) \stackrel{S_1}{\bowtie} (\{S1, S4\} - \{S1\}) \\ = (S_1, S_2) \stackrel{S_1}{\bowtie} (\{S4\}) \\ = \{(S4, P2), (S4, P4), (S4, P5)\}.$

5. Because there is only one attribute Suppliers.sname in $T(Q)$, we eliminate unnecessary join operations and perform

$$\pi_{sname}(\subset S_1, S_2 \supset \stackrel{sno}{\bowtie} RE_1) \\ = \pi_{sname}(\{((S4, P2), (S4, P4), (S4, P5))\} \stackrel{sno}{\bowtie} \\ \sigma_{city=London}(Suppliers)) \\ = \{Clark\}. \quad \square$$

5 Conclusions and Future Work

We study the inter-relationship between the natural language constructs of a query and the E-R conceptual schema. The basic constructs of English sentences can be mapped onto E-R schemas in a natural way. We develop a logical form by extending the E-R concepts to capture natural language semantics and describe a processing model for the query transformation process. In this processing model, when the target database is changed we only have to change the dictionary and the ER/SRF. The front-end parser/mapper and the query transformer remain unchanged.

English sentences may also be mapped onto the universal relation [16], in which the entire database is imagined to be kept in a single relation. But this needs to further transform the universal query command into the actual stored schema. Moreover, other intermediate forms suffer from the bias to natural language constructs and much effort is needed to transform them into database query languages. In comparison, our logical form has the merit that it is represented in a form similar to the ERD and can be efficiently transformed into the relational algebra.

Finally, an extension for mapping the natural language constructs of a query onto the schema generated by an object-oriented design methodology (e.g. the one proposed by Blaha, et al. [2]) will be investigated in the near future. Besides, by combining the framework presented by Zvieli and Chen [20], our work can be extended to process a natural language query involving a modifier like 'almost', 'very', or 'nearly'. This combination is served as a step toward analyzing the use of modifiers, which are fuzzy in natural, to communicate with fuzzy databases.

Acknowledgement

The authors wish to thank the anonymous reviewers for their valuable comments.

References

- [1] P.A. Bernstein and D.M.W. Chiu, Using Semi-Joins to Solve Relational Queries, *Journal of the ACM* 28 (1) (1981) 25-40.
- [2] M.R. Blaha, et al., Relational Database Design Using an Object-Oriented Methodology, *Comm. ACM* 31 (4) (1988) 427.
- [3] P.P. Chen, The Entity-Relationship Model — Toward a Unified View of Data, *ACM TODS* 1 (1) (1976) 9-36.
- [4] P.P. Chen, English Sentence Structure and E-R Diagrams, *Information Sciences* 29 (2) (1983) 127-149.
- [5] A.L.P. Chen, et al., Distributed Query Processing in a Multiple Database System, *IEEE Journal on Selected Areas in Communications* 7 (3) (1989) 390-398.
- [6] F.J. Damerou, Problems and Some Solutions in Customization of Natural Language Database Front Ends, *ACM Trans. Office Information Systems* 3 (2) (1985) 165-184.
- [7] C.J. Date, *A Guide to INGRES* (Addison-Wesley, MA, 1987).
- [8] C.J. Date, *An Introduction to Database Systems* (Addison-Wesley, MA, 5th ed., 1990).
- [9] B.J. Grosz, et al., TEAM: An Experiment in the Design of Transportable Natural-Language Interfaces, *Artificial Intelligence* 32 (1) (1987) 173-243.
- [10] G. Jakobson et al., An Intelligent Database Assistant, *IEEE Expert* 1 (2) (1986) 65-79.
- [11] R.A. Kowalski, *Logic for Problem Solving* (North-Holland, Amsterdam, 1979).
- [12] M.C. McCord, Using Slots and Modifiers in Logic Grammars for Natural Language, *Artificial Intelligence* 18 (3) (1982) 327-367.
- [13] F.C.N. Pereira and D.H.D. Warren, Definite Clause Grammars for Language Analysis — A Survey of the Formalism and a Comparison with Augmented Transition Networks, *Artificial Intelligence* 13 (3) (1980) 231-278.
- [14] G. Piatetsky-Shapiro and G. Jakobson, An Intermediate Database Language and Its Rule-based Transformation to Different Database Languages, *Data & Knowledge Engineering* 2 (1987) 1-29.
- [15] F.S.C. Tseng, S.Y. Lee and W.P. Yang, DELICIOUS: An Intermediate Code Scheme for Heterogeneous Database Systems, *Proc. International Computer Symposium, Taiwan, ROC* (1988) 998-1003.
- [16] J.D. Ullman, *Principles of Database and Knowledge-Base Systems* (Computer Science Press, Vol. 2, Rockville, MD, 1988).
- [17] M. Wallace, *Communicating with Databases in Natural Language* (Ellis Horwood, England, 1984).
- [18] P.H. Winston, *Artificial Intelligence* (Addison-Wesley, MA, 1984).
- [19] W.A. Wood, Transition Network Grammars for Natural Language Analysis, *Comm. ACM* 13 (10) (1970) 591-606.
- [20] A. Zvieli and P.P. Chen, Entity-Relationship Modeling and Fuzzy Databases, *Proc. IEEE Int. Conf. Data Engineering* (1986) 320-327.