**Association for Information Systems**
**AIS Electronic Library (AISeL)**

12-31-1994

# The Query Cube: A Framework for Assessing User Productivity with Database Information Retrieval

Hock Chan
*National University of Singapore*

Bernard Tan
*University of South Africa*

M. Grobler
*University of Pretoria*

Jonathan Miller
*University of Cape Town*

Mayuri Odedra-Straub
*Odedra-Straub Research and Consultancy*

*See next page for additional authors*

**Authors**

Hock Chan, Bernard Tan, M. Grobler, Jonathan Miller, Mayuri Odedra-Straub, and Jackie Phahlamohlaka

# THE QUERY CUBE: A FRAMEWORK FOR ASSESSING USER PRODUCTIVITY WITH DATABASE INFORMATION RETRIEVAL

Hock C. Chan
Bernard C. Y. Tan
Kwok-Kee Wei
Department of Information Systems and Computer Science
National University of Singapore

## ABSTRACT

Three key factors that affect user productivity on database information retrieval are representation realism, expressive ease, and task complexity. Representation realism is the level of abstraction used in formulating queries. Expressive ease is the syntactic flexibility of a query language. Task complexity is the level of difficulty of queries. These factors formed a three dimensional query cube. A laboratory experiment was conducted to evaluate user productivity on database information retrieval corresponding to different vertices of the query cube. The results show that the query cube is a viable framework for assessing user productivity, both on effectiveness and efficiency perspective.

## 1. INTRODUCTION

The growing use of computers by non-technical people makes it critical to provide user-friendly interfaces (Codd 1974; Cuff 1980; Mayer 1981; Moran 1981). For databases, a good user-database interface allows fast and accurate query formulation. For productive interaction, users require data knowledge and language knowledge (Ogden 1986). A major obstacle for end-users is the lack of required knowledge.

Data knowledge is provided by the data model presented to users. The de facto model is the relational model. However, several researchers have noted the difficulty of the relational model in capturing real world semantics (Kent 1979). The need for semantic models has seen the introduction of the entity-relationship (ER) model (Chen 1976) which has received widespread interest. Claims of its superiority over the relational model have been supported in some studies. Batra, Hoffer and Bostrom (1990) and Jarvenpaa and Machesky (1986) found that ER model users designed better than relational model users. Chan, Wei and Siau (1993) reported higher retrieval productivity for ER model users compared to relational model users. However, Jih et al. (1989) found no difference in productivity between users of these models.

Besides the data model, query language also affects user productivity. Linear keyword languages (LKL), e.g., SQL, have been used for most database retrieval operations. An LKL has a restricted syntax with a limited set of keywords. Users encounter difficulties using SQL because of its syntactic rigidity (Barbary 1987). The need to learn a new syntax lowers user productivity (Malhotra and Wladawsky 1975). An alternative is a natural language (NL) interface that allows novice users to formulate database queries with minimal training.

NL interfaces are becoming more common. However, human factor studies on NL database queries have produced mixed results. Some supported the use of NL (Suh and Jenkins 1992; Vassiliou et al. 1983) while others supported the use of LKL (Jarke et al. 1985; Shneiderman 1980; Small and Weldon 1983). Moreover, some reported no significant differences (Turner et al. 1984). These conflicting results raise an interesting research question: Are there advantages in using NL?

This study examines the impact of representation realism (ER versus relational model), expressive ease (NL versus LKL), and task complexity (simple versus complex queries) on user productivity (query accuracy and time taken). With more NL interfaces becoming available for relational databases, it is important to investigate whether data models of higher realism are still useful. These three factors are combined into a query cube in section 2. Section 3 illustrates the research hypotheses and methodology. The results are analyzed in section 4 and discussed in section 5. Section 6 concludes this paper.
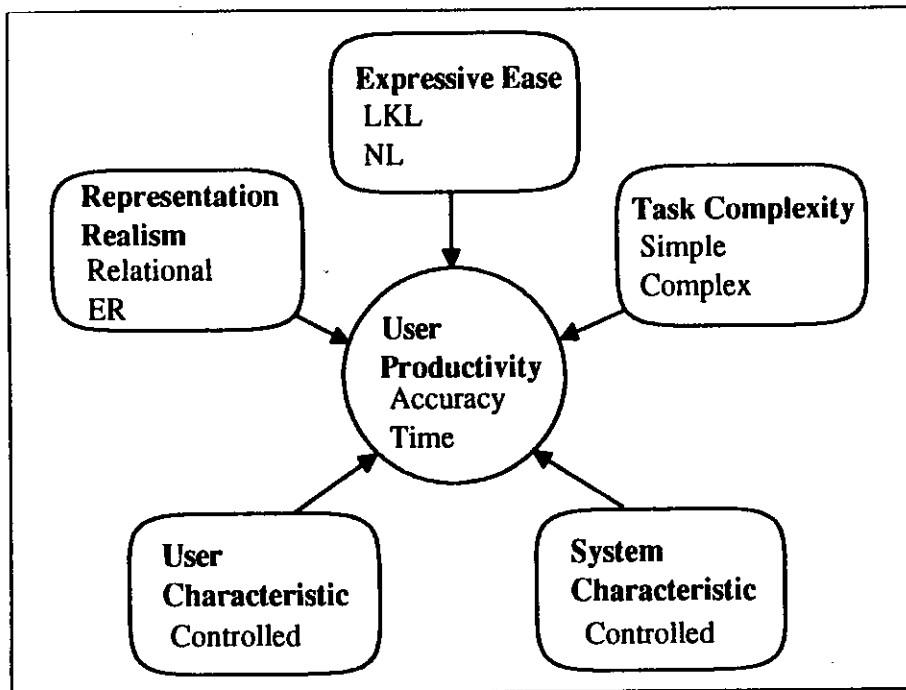
**Figure 1. Research Model**

## 2. LITERATURE REVIEW

Studies on the use of query languages have spanned two decades. These studies have measured user productivity by query accuracy and time taken. The important factors studied are data model (including query languages), task, user characteristics and system characteristics (Chan, Wei and Siau 1993; Reisner 1981). The present study is on data model and task. The other two factors are left for further studies.

The research model showing the three factors of representation realism, expressive ease, and task complexity is shown in Figure 1. The following subsections provide detailed definitions of these factors and review existing findings on their effects on user productivity.

### 2.1 Representation Realism

For databases, representation realism is measured by the level of abstraction of the interaction between the user and the database system (Chan, Wei and Siau 1993). It indicates how close a representation is to the concept being represented. The three main levels of abstraction, in ascending order of realism, are the physical, logical and conceptual levels.

The ER model belongs to the conceptual level and has a higher degree of representation realism; the relational model

is associated with the logical level and has a lower degree of representation realism (Batini, Ceri and Navathe 1992; Hughes 1991; Vossen 1991). The few experiments comparing the effects of representation realism for database retrieval operations have not produced consistent results (Chan, Wei and Siau 1993; Jih et al. 1989). Further research is needed.

### 2.2 Expressive Ease

Expressive ease is the amount of syntactic freedom accorded to users by a query language. NL is high in expressive ease because it has flexible syntactic rules and users have years of experience expressing themselves using NL. LKL is low in expressive ease because it has a very rigid set of syntactical rules and users may be unfamiliar with it.

The experimental results on NL versus LKL have been inconsistent. Shneiderman (1980), Small and Weldon (1983), Turner et al. (1984), and Vassiliou et al. (1983) reported no differences in the number of correct queries between English and SQL. These studies highlighted the importance of user knowledge on the application domain, especially for NL users. Many irrelevant queries were posted due to a lack of domain knowledge.

Jarke et al. found that SQL subjects were more than twice as successful in accomplishing their tasks as NL subjects.

This could be partially attributed to system rather than language differences. NL subjects encountered more frequent system problems than SQL subjects. In contrast, Suh and Jenkins reported better results for NL users. SQL was compared with a commercial NL interface (R:Base Clout) for data retrieval with novice users. Queries were wide-ranging, covering simple conditions, built-in functions, and/or conditions, and combination queries.

A possible explanation for the inconsistent results is the biases against NL in the early experiments. These biases include the use of prototype NL systems with many system problems, the lack of database training, and the many restrictions on NL syntax. Moreover, some experiments had small sample sizes resulting in low statistical power (Cohen 1988). Therefore, these results must be interpreted cautiously.

It is important to distinguish the two totally different methods by which NL can enhance user productivity. One method is by allowing syntactic freedom. The other is by raising the representational realism. Through customization in NL interfaces, database professionals can hide database operations. For example, a user may be able to issue the command *find the suppliers of Royal tires* without specifying any relations or joins. In this study, these two methods are separated. The NL subjects are required to express elements of the data model. That is, relational subjects must specify the relations and joins, while ER subjects must specify the entities and relationships.

## 2.3 Task Complexity

Reisner (1977) classified SQL queries into eleven types and recommended that SQL be treated as a layered language. The first (simple) layer includes simple mapping of returning data values when a known data value for another field is supplied and Boolean queries permitting and/or connectives. The second (complex) layer includes combination queries which use the output of one query as input for another, groupings which collect items with a common domain value, and universal quantifications which correspond to the "for all" concept of first order calculus. The complexity comes from the use of operations such as nesting and grouping that are assessed to be difficult. This classification of simple and complex queries has been validated in many studies on SQL (Welty 1985; Welty and Stemple 1981). This two-layer classification of query difficulty is not totally satisfactory. Queries vary from the very simple to the very difficult. Much work needs to be done to provide some quantitative measure of query difficulty.

Vassiliou et al. compared SQL with USL, a prototype NL system. They found that SQL was more suited to complex queries while USL was better for simple queries. Small and Weldon (1983) reported that, for "simple mapping" and "composition" problems, queries using SQL took less time to formulate. Suh and Jenkins found that NL subjects were faster and more accurate than SQL subjects for simple, and/or, and composition queries. No differences were found for queries with functions.

## 2.4 The Query Cube

Representation realism, expressive ease, and task complexity can be combined into a three-dimensional query cube, as in Figure 2. The representation realism dimension depicts the level of abstraction permitted when formulating database queries. High representation realism is characterized by the ability to express queries using concepts from the user's world rather than data structures from the database system. The expressive ease dimension illustrates the syntactic freedom allowed in expressing database queries. High expressive ease results from the use of a familiar language and the presence of syntax flexibility. Low expressive ease results from rigid syntax and limited keywords. It is also possible to view representational realism as the semantics and expression ease as the syntactic components of a query. The task complexity dimension measures the level of difficulty of database queries. Queries can be placed along this dimension based on Reisner's (1977) classification scheme.

Any database query can be represented with a point on the query cube. User productivity is expected to be high when representation realism and expressive ease are high and when task complexity is low. The distance between the point and the origin of the query cube provides an indication of user productivity. The conflicting results on representation realism and expressive ease and the lack of research on task complexity suggest that these dimensions are worthy of further research. Also, it should be noted that the relative impact of the three dimensions on user productivity is uncertain. Further research is needed to shed light on this issue.

## 3. RESEARCH HYPOTHESES AND METHODOLOGY

The three independent variables investigated in this study are representation realism, expressive ease, and task complexity. Representation realism is varied at two levels: high (ER model) and low (relational model). Expressive ease is examined at two levels: high (NL) and low (LKL).
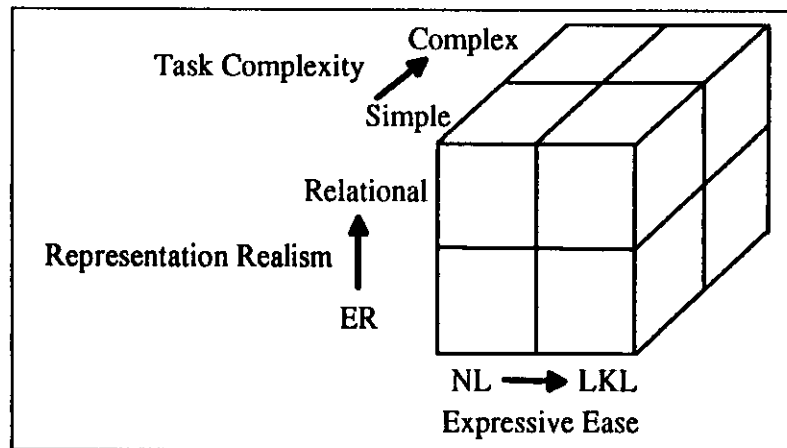
**Figure 2. The Query Cube**

Task complexity is also studied at two levels: simple and complex (based on Reisner's [1977] classification scheme). The two dependent variables are measures of user productivity: query accuracy and time taken. These two measures gauge the effectiveness and efficiency of users in formulating database information retrieval operations respectively. The controlled variables are human and database system characteristics.

## 3.1 Research Hypotheses

The hypotheses are based on the proposition that the closer the query is to the origin in the query cube, the easier it will be for users to write the query. Ease of query writing is measured by the accuracy of the query and the time taken to write the query.

H1a:   Users using the ER model will achieve greater query accuracy than users using the relational model.
H1b:   Users using NL will achieve greater query accuracy than users using LKL.
H1c:   Users performing simple queries will achieve greater query accuracy than users performing complex queries.

H2a:   Users using the ER model will take less time than users using the relational model.
H2b:   Users using NL will take less time than users using LKL.
H2c:   Users performing simple queries will take less time than users performing complex queries.

## 3.2 Experiment Design

The research design was a 2x2x2 factorial controlled laboratory experiment, with task complexity being a repeated measure. Table 1 shows the number of subjects in each treatment.

## 3.3 Independent Variables

The ER model corresponds to a higher (conceptual) level of abstraction while the relational model is associated with a lower (physical) level of abstraction (Chan, Wei and Siau 1993). Hence, high and low representation realism are operationalized using the ER model and relational model respectively.

With the relational model, low expressive ease was set using SQL, the standard language for relational databases. A relational NL was used to permit high expressive ease. With the ER model, low expressive ease was set using knowledge query language (KQL) (Chan, Wei and Siau 1993), the equivalent of SQL for the ER model. High expressive ease was operationalized using an NL equivalent in syntax to conversational English and able to express concepts in the ER model.

Queries were classified into simple or complex based on Reisner's (1977) classification scheme. Simple queries were those included in the first layer of this scheme while complex queries were those corresponding to the second layer.

**Table 1. Experiment Design**

| Representation realism | Expressive ease | |
|---|---|---|
| | NL | LKL |
| ER model | 27 subjects | 30 subjects |
| Relational model | 29 subjects | 26 subjects |

## 3.4 Dependent Variables

Many empirical studies on database query languages have measured query accuracy as a dependent variable (Jarke et al. 1985; Suh and Jenkins 1992; Reisner 1977; Turner et al. 1984; Welty 1985). Query accuracy reflects user effectiveness in formulating queries. It was assessed using the Welty correctness scheme of grading. Both semantic and syntactic accuracy were considered during the grading process. Each query was given a score ranging from 0 to 5, with a higher score indicating better query accuracy. The scores for simple and complex queries for each subject were obtained by averaging the scores of all his simple and complex queries respectively. An author and an independent grader determined the accuracy of all queries. The correlation between the two sets of grades for accuracy were very high. Hence, the scores given by the independent grader were adopted.

Besides query accuracy, the time taken to formulate queries has also been examined in many empirical studies on database query languages (Chan, Wei and Siau 1993; Jarke et al. 1985; Small and Weldon 1983; Vassiliou et al. 1983). This measure assessed user efficiency in expressing queries. The time taken to formulate a query was measured in seconds and recorded automatically by the database system. Timing started when the question was displayed on the screen and ended when the subject typed ALT_A to signal the completion of his query. The mean times taken for simple and complex queries, for each subject, were computed by averaging the time taken for all his simple and complex queries respectively.

## 3.5 Controlled Variables

The important controlled variables were human and database system characteristics. Human factors were controlled by randomly assigning subjects to different treatments. A total of 112 first-year computer science undergraduates participated as subjects in this study. A large sample size was used to ensure adequate statistical power (Cohen 1988). The subjects had little knowledge of databases prior

to this study. They were given course credit for their participation. They were informed that the credits would depend on their performance in the test. Accuracy would be the main factor and time would be a minor factor.

Database system characteristics were kept constant by using a common database query environment for all treatments. The database system adopted a user-friendly Windows interface to ensure that subject performance would not be hindered by system limitations. System limitations had been a confounding factor contributing to conflicting results in earlier empirical studies.

## 3.6 Experiment Procedure

Four experiment sessions, each corresponding to a different treatment shown in Table 1, were conducted. Each session began with a short training period. Subjects were given a fifteen minute briefing on the important concepts of the data model (ER or relational model). This was followed by a forty-five minute practice period where subjects formulated queries on the database system using the language corresponding to their respective treatments. A sample database schema (ER diagram or relational schema) together with five simple and five complex queries were provided to each subject for this purpose. Besides providing database knowledge to subjects, the sessions also familiarized them with the query environment.

After the training, subjects proceeded to perform the actual queries. Each subject was given a new database schema (ER diagram or relational schema) together with a description of the database domain pertaining to the actual queries. Each subject had to formulate six simple followed by six complex queries. These queries were presented in the same order for all treatments. No attempt was made to control carry over effects by varying the sequence of simple and complex queries. Instead, all subjects performed simple followed by complex queries because this reflects closely the usage patterns of real-life database users. Subjects performed one query at a time by entering the query into the system. They could refer to the training materials.

65

They were also given paper and pencils. The end of a query was signalled using the ALT_A key. Upon completion of a query, subjects proceeded to the next query. No time limit was imposed for a query.

## 4. DATA ANALYSES

All statistical tests were performed at a 5% level of significance. An overall MANOVA test detected significant main effects due to representation realism (F = 202.55, p = 0.0001), expressive ease (F = 9.60, p = 0.0001), and task complexity (F = 90.79, p = 0.0001). This test also found significant two-way interactions involving representation realism and expressive ease (F = 20.33, p = 0.0001) and representation realism and task complexity (F = 37.35, p = 0.0001) as well as a significant three-way interaction involving all three independent variables (F = 3.21, p = 0.0421). These significant results permitted separate ANOVA tests to be carried out on query accuracy and time taken. No transformations (Weisberg 1985) could be found to allow query accuracy and time taken to simultaneously meet the homogeneity and normality requirements (Neter, Wasserman and Kutner 1990) of the ANOVA test. Hence, all significant results detected by the ANOVA test were confirmed using the non-parametric Mann-Whitney test (Siegel and Castellan 1988). The means and standard deviations of the dependent variables for the different treatments are shown in Table 2. Results of the ANOVA test are summarized in Table 3.

The score for query accuracy ranges from 0 to 5 with a higher score indicating better query accuracy. All three independent variables had significant main effects for query accuracy (see Table 3). The Mann-Whitney test confirmed the main effects due to representation realism (Chi-square = 45.20, p = 0.0001), expressive ease (Chi-square = 4.96, p = 0.0259), and task complexity (Chi-square = 33.05, p = 0.0001). Higher query accuracy was obtained in conditions of high representation realism, high expressive ease, or low task complexity. Hypotheses 1a, 1b, and 1c were all supported.

Besides the main effects, there were two way interactions between representation realism and expressive ease, representation realism and task complexity, and expressive ease and task complexity (see Table 3). These interactions were analyzed using the method of analysis of simple effects (Keppel 1991). Results of these analyses are summarized in Table 4. All significant results were confirmed by the Mann-Whitney test. These findings showed that (1) representation realism affected query accuracy under all conditions of expressive ease and task complexity, (2) task complexity had an impact on query accuracy under all conditions of representation realism and expressive ease, and (3) expressive ease affected query accuracy only when representation realism was high or when task complexity was high. Among the three independent variables, expressive ease seemed to have the least widespread impact on query accuracy.

### Table 2. Means and Standard Deviations of Dependent Variables

| Representation realism | Expressive ease | Task complexity | Query accuracy | Time taken |
|---|---|---|---|---|
| ER Model | NL | Simple | 4.65 (0.47) | 65 (24) |
| | | Complex | 4.00 (0.73) | 70 (22) |
| | LKL | Simple | 3.83 (1.07) | 112 (28) |
| | | Complex | 3.12 (1.18) | 134 (67) |
| Relational Model | NL | Simple | 3.30 (0.82) | 158 (52) |
| | | Complex | 2.18 (1.50) | 294 (74) |
| | LKL | Simple | 3.80 (0.72) | 156 (50) |
| | | Complex | 1.58 (1.27) | 256 (73) |

Table 3. Results of the ANOVA Test

| Factor | Query accuracy | | Time taken | |
|---|---|---|---|---|
| | F-value | p-value | F-value | p-value |
| Representation realism (RR) | 74.16 | 0.0001** | 290.77 | 0.0001** |
| Expressive ease (EE) | 10.74 | 0.0012** | 6.25 | 0.0132* |
| Task complexity (TC) | 72.76 | 0.0001** | 87.40 | 0.0001** |
| RR x EE | 8.26 | 0.0045** | 28.24 | 0.0001** |
| RR x TC | 12.57 | 0.0005** | 54.91 | 0.0001** |
| EE x TC | 4.58 | 0.0335* | 0.38 | 0.5357 |
| RR x EE x TC | 3.66 | 0.0571 | 3.59 | 0.0594 |

*$p < 0.05$        **$p < 0.01$

Table 4. Results of Interaction Analyses on Query Accuracy

| Interaction examined | Variable controlled | Representation realism | Expressive ease | Task complexity |
|---|---|---|---|---|
| RR x EE | RR<br>  ER<br>  Relational | <br>NA<br>NA | <br>p = 0.0001**<br>p = 0.8378 | <br>NA<br>NA |
| | EE<br>  NL<br>  LKL | <br>p = 0.0001**<br>p = 0.0024** | <br>NA<br>NA | <br>NA<br>NA |
| RR x TC | RR<br>  ER<br>  Relational | <br>NA<br>NA | <br>NA<br>NA | <br>p = 0.0004**<br>p = 0.0001** |
| | TC<br>  Simple<br>  Complex | <br>p = 0.0001**<br>p = 0.0001** | <br>NA<br>NA | <br>NA<br>NA |
| EE x TC | EE<br>  NL<br>  LKL | <br>NA<br>NA | <br>NA<br>NA | <br>p = 0.0003**<br>p = 0.0001** |
| | TC<br>  Simple<br>  Complex | <br>NA<br>NA | <br>p = 0.4608<br>p = 0.0192* | <br>NA<br>NA |

*$p < 0.05$        **$p < 0.01$
NA: Not applicable for the interaction examined

**Table 5. Results of Interaction Analyses on Time Taken**

| Interaction examined | Variable controlled | Representation realism | Expressive ease | Task complexity |
|---|---|---|---|---|
| RR x EE | RR<br>　ER<br>　Relational | <br>NA<br>NA | <br>p = 0.0001**<br>p = 0.2350 | <br>NA<br>NA |
|  | EE<br>　NL<br>　LKL | <br>p = 0.0001**<br>p = 0.0001** | <br>NA<br>NA | <br>NA<br>NA |
| RR x TC | RR<br>　ER<br>　Relational | <br>NA<br>NA | <br>NA<br>NA | <br>p = 0.1248<br>p = 0.0001** |
|  | TC<br>　Simple<br>　Complex | <br>p = 0.0001**<br>p = 0.0001** | <br>NA<br>NA | <br>NA<br>NA |
| EE x TC | EE<br>　NL<br>　LKL | <br>NA<br>NA | <br>NA<br>NA | <br>NT<br>NT |
|  | TC<br>　Simple<br>　Complex | <br>NA<br>NA | <br>NT<br>NT | <br>NA<br>NA |

*p < 0.05　　**p < 0.01
NA = Not applicable for the interaction examined
NT = Not tested because no interaction was detected

Time taken is measured in seconds, with a lower value indicating faster query formulation. All three independent variables had significant main effects for time taken (see Table 3). The Mann-Whitney test confirmed the main effects due to representation realism (Chi-square = 110.80, p = 0.0001), expressive ease (Chi-square = 7.49, p = 0.0062), and task complexity (Chi-square = 16.92, p = 0.0001). Queries were formulated faster when representation realism was high, expressive ease was high, or when task complexity was low. Hypotheses 2a, 2b, and 2c were all supported.

Besides the main effects, there were two way interactions due to representation realism and expressive ease and representation realism and task complexity (see Table 3). The method of analysis of simple effects (Keppel 1991) was also used to analyze these interactions. Results of these analyses are summarized in Table 5. The Mann-Whitney test confirmed all significant results. These findings showed that (1) representation realism affected time taken under all conditions of expressive ease and task complexity, (2) expressive ease had an impact on time taken only when representation realism was high, and (3) task complexity affected time taken only when representation realism was low. Among the three independent variables, representation realism seemed to have the most widespread impact on time taken.

## 5. DISCUSSION AND IMPLICATIONS

### 5.1 Results of Main Analyses

The results of this experiment confirm that representation realism, expressive ease, and task complexity are key factors influencing user productivity with database information retrieval, from both an effectiveness (query accuracy) and an efficiency (time taken) perspective. Users of the ER model achieved better productivity than user of the relational model. Users of NL attained higher productivity than users of LKL. Users performing simple queries also experienced greater productivity than users formulating complex queries.

68

With the relational model, users had to logically specify the joins by linking the relations based on common values. This was rather procedural and required the construction of a complicated expression to perform a conceptually simple task. Users of the ER model need not perform these cumbersome procedures, thereby resulting in significantly better query accuracy and shorter query time. These findings on the importance of representation realism confirm those of Chan, Wei and Siau. The finding that ER model users are faster then relational users is also consistent with Jih et al. Based on these finding, it is recommended that user productivity could be enhanced by adopting a higher level of representation realism. One way of attaining this is to use an ER model instead of a relational model for the user-database interface.

The significant findings on expressive ease could be due to inherent advantages that accompany NL. With NL, there is no need to learn a new syntax. Moreover, there is syntax flexibility in the sense that users need not be constrained by a very rigid set of syntax rules. Common deviations from formally correct grammar could be tolerated as long as the intent of the request was clear. Another contributing factor could be that NL permits the use of synonyms. In contrast, LKL has a very specific set of syntax rules where deviations are not allowed. These constraints could have contributed to a reduction of user productivity when formulating queries. These findings add to the small body of literature (Suh and Jenkins 1992) supporting the use of NL over LKL. Hence, user productivity could be enhanced by raising the level of expressive ease. More development effort toward the direction of NL could result in tangible benefits.

Higher query accuracy and shorter query time resulted with simple queries, compared to complex queries. Simple queries have fewer operations and require less knowledge and effort to formulate. Hence, users attained greater productivity with simple queries. These results are consistent with the findings of Jih et al. and of Suh and Jenkins on the effect of task complexity.

## 5.2 Results of Interaction Analyses

In terms of query accuracy, higher representation realism and lower task complexity contributed to enhanced user performance under all conditions. However, expressive ease had no impact when representation realism was low or when task complexity was low. When representation realism is low, many data operations are required for formulating queries resulting in a greater possibility of introducing errors. These additional required operations might have reduced whatever gains come with the use of

NL. When queries are simple, users of LKL are generally able to attain a high level of query accuracy. Hence, additional benefits that accrue from the use of NL are likely to be minimal.

In terms of time taken, higher representation realism raised user performance under all circumstances. However, expressive ease had no impact when representation realism was low and task complexity had no impact when representation realism was high. A low level of representation realism necessitates numerous data operations requiring much effort during query formulation. Savings in time brought about through the use of NL might have been offset by the need to perform many data operations. A high level of representation realism permits users to skip many cumbersome procedural operations. Users are generally able to formulate queries using a few simple operations for both simple and complex queries. Few benefits result from the use of simple queries.

## 5.3 Implications for the Query Cube

The three dimensions of the query cube are good predictors of user productivity with database information retrieval. Users are more productive when representation realism is high, when expressive ease is high, or when task complexity is low. Therefore, the origin of the query cube represents a point where user productivity is optimal. Departure from the origin along any of the three dimensions results in a deterioration of user productivity. Any query can be represented as a point in the query cube based on these three dimensions. The distance between this point and the origin of the query cube provides a preliminary indication on user productivity for that particular query situation.

Besides serving as a framework to predict user productivity, the query cube provides a direction for future research and development on database queries. Task complexity is a factor that cannot be controlled. Complex queries could not be substituted with simple queries to raise user productivity. However, user productivity could be raised by employing a higher level of representation realism and expressive ease. Research and development effort aimed at enhancing representation realism, perhaps by improving the ER model or perhaps by introducing newer and better models, is likely to be worthwhile in the long term. Likewise, although a truly NL interface does not exist at present, research and development work aimed at the creation of query languages resembling the English language is likely to pay off in the long term. Expressive ease had a major impact on user productivity when representation realism was high. Hence, users could best enhance their

productivity by adopting a database model with high representation realism, such as the ER model, before attempting to use a database language with high expressive ease, such as NL. Nevertheless, the axis of the query cube corresponding to high representation realism and high expressive ease provides a target that researchers and users could approach.


## 6. CONCLUSION

Results of this study confirm the findings from prior studies that representation realism, expressive ease, and task complexity are important individual factors determining user productivity with database information retrieval. Moreover, the results of this study contribute additional knowledge on how these individual factors interact to influence user productivity. For instance, expressive ease could only affect query accuracy and time taken under certain combinations of representation realism and task complexity conditions.

This study proposes a three dimensional query cube as a framework to predict user productivity with database information retrieval. The predictive potential of the query cube was supported by empirical data. However, predictions made using the query cube should be considered tentative, rather than definite, and interpreted with caution. Although the query cube highlights three important factors in predicting user productivity, the relative importance of these factors is uncertain and is a topic requiring further research.

Future studies could also assess the query cube to see if it is able to predict other aspects of user performance besides productivity. Examples of other performance variables are user confidence during query formulation, user ability to recover from mistakes when prompted with error messages, and user inclination to learn from their experience with prior queries. Given that user productivity dropped on departure from the origin of the query cube, an error analysis could be done to see which kinds of errors occurred at each vertex of the query cube. Such an analysis would help vendors of database systems better predict the kind of difficulty their users are likely to encounter so that relevant support, in the form of help facilities or user manuals, could be developed. With an increasing number of end-users performing database queries (Ahrens and Sankar 1993; Cronan and Douglas 1990), it becomes critical that more research effort be directed at understanding factors affecting user productivity and seeking ways to improve it. This study is a step in that direction.

## 7. REFERENCES

Ahrens, J. D., and Sankar, C. S. "Tailoring Database Training for End Users." *MIS Quarterly*, Volume 17, Number 4, 1993, pp. 419-439.

Barbary, C. "A Database Primer on Natural Language." *Journal of Systems Management*, Volume 38, Number 4, 1987, pp. 20-25.

Batini, C.; Ceri, S.; and Navathe, S. B. *Conceptual Database Design, An Entity Relationship Approach*. Menlo Park, California: Benjamin Cummings Publishing Company Inc., 1992.

Batra, D.; Hoffer, J.; and Bostrom, R. "Comparing Representations with Relational and EER Models." *Communications of the ACM*, Volume 33, Number 2, 1990, pp. 126-138.

Chan, H. C.; Wei, K. K.; and Siau, K. L. "User-Database Interface: The Effect of Abstraction Levels on Query Performance." *MIS Quarterly*, Volume 17, Number 4, 1993, pp. 441-464.

Chen, P. P. "The Entity-Relationship Model: Toward a Unified View of Data." *ACM Transactions on Database Systems*, Volume 1, Number 1, 1976, pp. 9-36.

Codd, E. F. "Seven Steps to Rendezvous with the Casual User." In J. W. Klimbie and K. L. Koffeman (Editors), *Data Base Management*. Amsterdam: North-Holland, 1974, pp. 179-199.

Cohen, J. *Statistical Power Analysis for the Behavioral Sciences*. Hillsdale, New Jersey: Lawrence Erlbaum Associates, 1988.

Cronan, T., and Douglas, D. "End User Training and Computing Effectiveness in Public Agencies: An Empirical Study." *Journal of MIS*, Volume 6, Number 4, 1990, pp. 21-40.

Cuff, R. N. "On Casual Users." *International Journal on Man-Machine Studies*, Volume 12, Number 2, 1980, pp. 163-187.

Hughes, J. G. *Object-Oriented Databases*. Englewood Cliffs, New Jersey: Prentice-Hall, 1991.

Jarke, M.; Turner, J. A.; Stohr, E. A.; Vassiliou, Y.; White, N. H.; and Michielsen, K. "A Field Evaluation of Natural

Language for Data Retrieval." *IEEE Transactions on Software Engineering*, Volume 11, Number 1, 1985, pp. 97-113.

Jarvenpaa, S. L., and Machesky, J. J. "End User Learning Behavior in Data Analysis and Data Modelling Tools." In L. Maggi, R. Zmud, and J. Wetherbe (Editors), *Proceedings of the Seventh Annual International Conference on Information Systems*, San Diego, California, 1986, pp. 152-167.

Jih, K.; Bradbard, D.; Snyder, C.; and Thompson, N. "The Effects of Relational and Entity-Relationship Data Models on Query Performance of End Users." *International Journal on Man-Machine Studies*, Volume 31, Number 3, 1989, pp. 257-267.

Kent, W. "Limitations of the Record Based Information Models." *ACM Transactions on Database Systems*, Volume 4, Number 1, 1979, pp. 107-131.

Keppel, G. *Design and Analysis: A Researcher's Handbook*. Englewood Cliffs, New Jersey: Prentice-Hall, 1991.

Malhotra, A., and Wladawsky, I. "The Utility of Natural Language Systems." Research Report #RE5739, IBM Technical Journal, Watson Research Center, New York, 1975.

Mayer, R. E. "The Psychology of Learning Computer Programming by Novices." *ACM Computing Surveys*, Volume 13, Number 1, 1981, pp. 121-141.

Moran, T. P. "An Applied Psychology of the User." *ACM Computing Surveys*, Volume 13, Number 1, 1981, pp. 1-11.

Neter, J.; Wasserman, W.; and Kutner, M. H. *Applied Linear Statistical Model: Regression, Analysis of Variance, and Experimental Designs*. Homewood, Illinois: Richard D. Irwin, 1990.

Ogden, W. C. "Implications of a Cognitive Model of Database Query: Comparison of a Natural Language, a Formal Language, and Direct Manipulation Interface." *ACM SIGCHI Bulletin*, Volume 18, Number 2, 1986, pp. 51-54.

Reisner, P. "Human Factors Studies of Database Query Languages: A Survey and Assessment." *Computing Surveys*, Volume 13, Number 1, March 1981, pp. 13-31.

Reisner, P. "Use of Psychological Experimentation as an Aid to Development of a Query Language." *IEEE Transactions on Software Engineering*, Volume 3, Number 2, 1977, pp. 218-229.

Shneiderman, L. *Software Psychology*. Cambridge, Massachusetts: Winthrop, 1980.

Siegel, S., and Castellan, N. J. *Nonparametric Statistics for the Behavioral Sciences*. New York: McGraw-Hill, 1988.

Small, J., and Weldon, L. "An Experimental Comparison of Natural and Structured Query Languages." *Human Factors*, Volume 25, Number 3, 1983, pp. 253-263.

Suh, K. S., and Jenkins, A. M. "A Comparison of Linear Keyword Language and Restricted Natural Language Data Base Interfaces for Novice Users." *Information Systems Research*, Volume 3, Number 3, 1992, pp. 252-272.

Turner, J. A.; Jarke, M.; Stohr, E. A.; Vassiliou, Y.; and White, N. H. "Using Restricted Natural Language for Data Retrieval: A Plan for Field Evaluation." In Y. Vassiliou (Editor), *Human Factors and Interactive Computer Systems*. Norwood, New Jersey: Ablex, 1984, pp. 163-190.

Vassiliou, Y.; Jarke M.; Stohr E. A.; Turner J. A.; and White N. H. "Natural Language for Database Queries: A Laboratory Study." *MIS Quarterly*, Volume 7, Number 4, 1983, pp. 44-61.

Vossen, G. *Data Models, Database Languages and DBMSs*. Reading, Massachusetts: Addison-Wesley, 1991.

Weisberg, S. *Applied Linear Regression*. New York: John Wiley & Sons, 1985.

Welty, C. "Correcting User Errors in SQL." *International Journal of Man-Machine Studies*, Volume 22, 1985, pp. 463-477.

Welty, C., and Stemple, D. W. "Human Factors Comparison of a Procedural and a Nonprocedural Query Language." *ACM Transactions on Database Systems*, Volume 6, Number 4, December 1981, pp. 626-649.