The Effects of State-Based and Event-Based Data Representation on User Performance in Query Formulation Tasks
Author(s): Gove N. Allen and Salvatore T. March
Source: *MIS Quarterly*, Vol. 30, No. 2 (Jun., 2006), pp. 269-290
Published by: Management Information Systems Research Center, University of Minnesota
Stable URL: http://www.jstor.org/stable/25148731
Accessed: 16-06-2015 12:10 UTC

# MIS Quarterly

# THE EFFECTS OF STATE-BASED AND EVENT-BASED DATA REPRESENTATION ON USER PERFORMANCE IN QUERY FORMULATION TASKS[1,2]

By:  Gove N. Allen
      A. B. Freeman School of Business
      Tulane University
      7 McAlister Drive
      New Orleans, LA  70118
      U.S.A.
      gallen@tulane.edu

      Salvatore T. March
      Owen Graduate School of Management
      Vanderbilt University
      Nashville, TN  37203
      U.S.A.
      sal.march@owen.vanderbilt.edu

## Abstract

*Ad hoc query formulation is an important task in effectively utilizing organizational data resources. To facilitate this task, managers and casual end-users are commonly presented with database views expressly constructed for their use. Differences in the way in which things, states, and events are represented in such views can affect a user's ability to understand the database, potentially leading to different levels of performance (i.e., accuracy, confidence, and prediction of the*

*accuracy of their queries). An experiment was conducted over the Internet involving 342 subjects from 6 universities in North America and Europe to investigate these effects. When presented with an event-based view, subjects expressing low or very low comfort levels in reading entity-relationship diagrams expressed confidence that better predicted query accuracy although there were no significant differences in actual query accuracy or level of confidence expressed.*

**Keywords:** Query formulation performance, event-based, state-based, artifact-based, data models, database user view, sense-making, E-R diagram

## Introduction

The ability to effectively utilize organizational data resources has become a major source of competitive advantage. Data warehouses, for example, commonly provide end-user access to organizational data in support of strategy formulation, real-time decision-making, and other management activities (Borthick et al. 2001; Gray and Watson 1998). Relational database management systems (RDBMS) provide the underlying technology for maintaining and accessing organizational data resources. The logical structure of such databases is often complex (Shasha 1996; Teorey et al. 1986). While there has been significant research on multidimensional and graphical interfaces to such databases (Speier and Morris 2003), SQL remains the standard language for specifying *ad hoc* queries. Accurately formulating queries in SQL is a challenging task (Chan et al. 1993; Leitheiser and March 1996; Siau et al. 2004) and the detrimental effects of using data from inappropriately formulated queries can be significant.

Relational database management systems facilitate SQL query formulation tasks by enabling the definition of *database views* or *virtual* tables (Halevy 2001) in addition to the definition of *base* or *implemented* tables. A *user view* is a collection of base and/or virtual tables that are visible to a user and against which the user specifies queries. The database management system automatically maps queries posed on user views into the base tables in the database. Hence, user views are logical-level constructs that provide different users with different conceptualizations of the same database. A user view defined for a specific set of users is frequently communicated to them by means of a conceptual-level diagram, for example, an entity-relationship (E-R) diagram (MicroStrategy 2003).

We study the effects of state-based and event-based user views and their corresponding conceptual-level diagrams on performance in query formulation tasks. A *state-based* view organizes data around things and the properties that define their states. An *event-based* view organizes data around events and the affected resources and agents.

Understanding the effects of different user views on query formulation performance is important to IS managers who must ensure that organizational data resources are properly used. We study three measures of query formulation performance: accuracy, confidence, and prediction of accuracy. The first two have been used in prior studies (e.g., Chan et al. 1993; Leitheiser and March 1996). The third, introduced in this study, indicates a user's proficiency at self-assessment and is a particularly important measure of performance for users who infrequently formulate queries against complex corporate data resources (Goodhue et al. 2000).

## Background ▉▉▉▉▉▉▉▉▉▉▉▉

Numerous data models have been posed to represent organizational data resources both during initial development and in subsequent use (Silberschatz and Korth 1996). The term *data model* has been used in at least three different ways in the literature. Often it is defined as a set of constructs and rules used to model a real-world domain at a specific level of abstraction (conceptual, logical, or physical). However, other definitions allow specific representations of a specific domain to be termed a *data model*; still others include recommendations about how to create the model. We use the term *data model* in its broadest sense and rely on specific terms as defined below when more precision is needed.

Using terminology presented by Wand and Weber (2002), a *conceptual-modeling grammar* defines the constructs and rules used to model a real-world domain (e.g., a business

application). The E-R model (Chen 1976), for example, is the basis for commonly used conceptual-modeling grammars focusing on data requirements (Antony and Batra 2002; Markowitz and Shoshani 1992). A *conceptual-modeling script* uses a conceptual-modeling grammar to represent a real-world domain. An E-R diagram, for example, is a conceptual-modeling script that uses the E-R grammar to represent the data requirements of a real-world domain.

Different E-R diagrams can correctly represent the data requirements of the same real-world domain (Bronts et al. 1995; Kent 1978). The E-R diagram produced for a real-world domain is determined, at least in part, by the conceptual-modeling method employed. A *conceptual modeling method* prescribes techniques for accomplishing a data modeling task, including procedures for identifying phenomena to be modeled and for mapping identified phenomena to a data modeling grammar's constructs.

We differentiate two types of conceptual-modeling methods. The first focuses on things and their descriptions (states), viewing the database as a snapshot of reality (Dey et al. 1995; Teorey et al. 1986). The second focuses on events and the affected resources and agents, viewing the database as a composite of transactions or economic events (McCarthy 1982). Both types of conceptual-modeling methods can utilize the E-R grammar which defines an entity as any thing or event "which can be distinctly identified" (Chen 1976, p. 10); however, E-R diagrams produced using them are typically quite different. Using the former type of method results in what we term a *state-based E-R diagram*. It contains an entity for each relevant thing in the real-world domain. Using the latter type results in what we term an *event-based E-R diagram*. It contains an entity for each relevant event and an entity for each of the affected things.

Consider, for example, a company that must keep track of building keys that have been assigned to its employees. Each key may be assigned to a particular employee and each employee may be assigned multiple keys. When an employee no longer needs a key, it is returned and may subsequently be assigned to a different employee. A state-based E-R method would identify "Key" and "Employee" as entities, representing the assignment of keys to employees as a relationship (Figure 1a). An event-based E-R method would additionally identify the events, "Assign Key" and "Return Key," and represent them as entities (Figure 1b).

The differences between these two types of conceptual-modeling methods are rooted in their different ontological foundations. Ontology has been used as a theoretical lens for examining various elements of information systems representations (e.g., Geerts and McCarthy 2002; Wand et al. 1999)
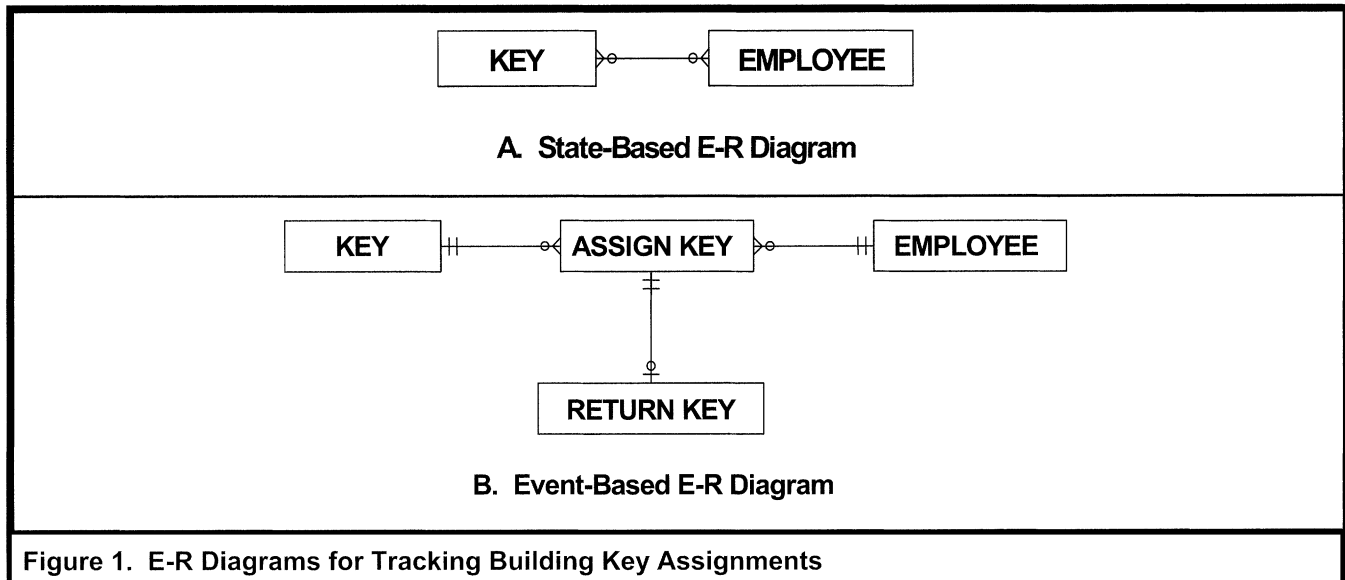
**A. State-Based E-R Diagram**

**B. Event-Based E-R Diagram**

**Figure 1. E-R Diagrams for Tracking Building Key Assignments**

and provides a basis for differentiating state-based and event-based E-R diagrams. An ontology defines a set of constructs used to represent real-world phenomena.

Ontologies commonly used in information system modeling include the concepts *thing* roughly corresponding to entity-instance and *property* roughly corresponding to attribute or relationship in the E-R grammar. The *state* of a thing is defined as the set of values of its properties at a point in time (Wand and Weber 1995). The concept of *event* is recognized in these ontologies; however, it is not consistently defined. The ontological works of Sowa (1999), Brody (1980), Tiles (1981), and Feibleman (1951) define things and events uniformly, allowing both to have existence (yielding identity) and properties. A **thing** *exists* at a given time, can be identified, and has properties. Similarly, an **event** *occurs* at a given time, can be identified, and has properties. From this ontological perspective events such as Assign Key and Return Key are appropriately represented as entities in an E-R diagram.

In contrast, Bunge's (1977) ontology and the information system ontology posed by Wand and Weber (1995) define an event as a "change of state of a thing" (p. 210) and conclude that event as an ontological construct is "not represented" in the E-R data modeling grammar (p. 217). Unlike things, which have existence (yielding the notion of *identity*) and properties, events themselves cannot have properties (Burton-Jones and Weber 1999; Wand et al. 1999). Adhering to such an ontology expressly precludes modeling events such as Assign Key and Return Key as entities in an E-R diagram. Debates about appropriate ontological underpinnings for

conceptual modeling have not yielded a basis for predicting human performance (Gemino and Wand 2005). To garner evidence with which to predict performance in query formulation we turn to the literature in psychology and human cognitive processing.

Humans have an innate competency for processing events. Human memory for events and past experiences is psychologically and physiologically different from human memory for facts and concepts (Nyberg 1998; Tulving 1983, 2002). Moreover, events are fundamental to narrative thinking (Robinson and Hawpe 1986) and to the representation of causality (Pillemer 1998; Ramesh and Browne 1999). Both are principal processes in human sense-making (Gee 1985). Furthermore, humans use this narrative or event processing competency as a powerful tool for verbal and written communication (Orr 1990).

Other human information processing competencies and limitations may also play important roles in determining query formulation performance. Two that have been considered in the information systems literature are construct overload and categorization. It can be argued that using entities to represent both things and events will result in construct overload and cause ambiguities in the model and deterioration of understanding and performance (Burton-Jones and Weber 1999; Wand and Weber 1995). This may be the case if people categorize events differently from things in their conceptualization of data, that is, if they do not ascribe existence or properties to both of them. Conversely if people conceptualize events as being *information bearing* (i.e., having

existence and properties), then this is not the case and an entity construct that treats them uniformly, as initially proposed by Chen (1976), is appropriate.

Information systems researchers have most often studied human categorization competencies in the context of classifying things (e.g., Parsons and Wand 2000). However, the human classification competency applies equally well to the categorization of human social interaction and the experience of events (Lakoff 1987). Thus, we expect this categorization competency to have a similar influence in reading categorized abstractions of a domain (such as E-R diagrams) whether they are state-based or event-based. Consequently we conjecture that if an event-based data model engages the human narrative competency, then it should result in improved understanding of a database and, therefore, improved query formulation performance.

## Prior Research ▮▮▮▮▮▮▮▮▮▮▮▮▮

Prior empirical research has investigated the effects of conceptual and logical models (grammars and scripts) on different types of database interactions including design, validation, understanding, and use in problem-solving and query formulation (e.g., Batra et al. 1990; Jih et al. 1989). Specifically, Chan et al. (1993) studied the effects of abstraction level on query performance. They found that subjects performed significantly better at *ad hoc* query formulation when interacting with a database at the conceptual level than at the logical level.

Kim and March (1995) studied the effects of two conceptual-modeling grammars, E-R and NIAM (Halpin 2001; Weber and Zhang 1991), on data modeling and validation tasks. They found that analysts using the E-R grammar produced models that were more accurate than analysts using the NIAM grammar. They found no significant performance differences between managers using models expressed in the E-R grammar and managers using models expressed in the NIAM grammar for validation tasks. The E-R and NIAM models used in this study were all state-based. This may be an explanation for the lack of significant results; users' understanding of a data model diagram (conceptual-modeling script) and the real-world domain it represents is more significantly affected by the ontological foundations of the conceptual-modeling *method* that produces it than by the conceptual-modeling *grammar* that expresses it.

Sinha and Vessey (1999) study end-user performance in developing conceptual-modeling scripts and corresponding logical-modeling scripts, comparing the E-R conceptual-modeling grammar and an object-oriented diagram (OOD) conceptual-modeling grammar with each other and with corresponding logical level grammars (relational and object-oriented text (OOT), respectively). They conclude that a conceptual-modeling grammar (E-R or OOD) results in superior modeling performance when compared to a logical-modeling grammar (relational or OOT). Performance is measured by the accuracy of the model produced using a given grammar for a set of specified constructs (e.g., entities/classes, attributes, and relationships). Accuracy is computed as a weighted percentage of correctly formulated instances of that construct in a subject's solution as compared to an expert's solution. We use a similar scheme to evaluate the accuracy of posed SQL queries.

More recently, researchers have studied differences in user performance that result from the use of ontologically diverse conceptual modeling methods. Bodart et al. (2001) and Gemino and Wand (2005) study the optional property construct in conceptual modeling. Conceptual-modeling methods conforming to Bunge's ontology preclude the use of optional properties while a number of commonly used conceptual modeling methods allow them. Both studies found that for problem solving tasks, precluding optional properties results in significantly better performance than allowing them. Bowen et al. (2004) studied the optional property construct in the context of query formulation performance. They found that for moderately complex models, precluding optional properties results in significantly worse performance than allowing them.

We similarly study the effects of ontological foundations on query formulation performance. However, rather than studying the effects of optional properties, we study the effects of differences in the ontological definition of the event construct.

## Research Methodology ▮▮▮▮▮▮▮▮▮▮

### *Research Model*

Our research model (Figure 2) is adapted from Chan et al. (1993). It asserts that query performance is influenced by characteristics of the data model and the user. Chan et al. study the abstraction level of a data model while we study its ontological foundations. Specifically, we study three ontological foundations: state-based, event-based, and a mix of the two that we term artifact-based. We use these to develop three distinct conceptual level (E-R) diagrams with corre-
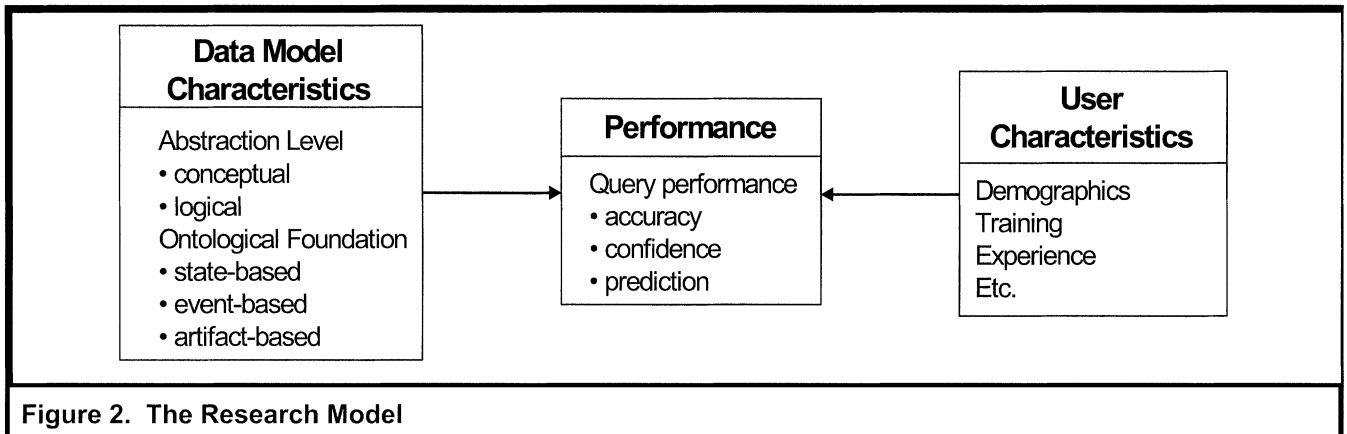
**Figure 2. The Research Model**

sponding logical level (relational) data models. Each pair comprises an experimental treatment.

We study query formulation performance using a single set of information requests (query requirements) for all subjects. User characteristics include demographics, training, experience, etc. We address differences in user characteristics by randomly assigning subjects to treatments.

## Dependent Variables: Query Performance Measures

Three variables measure query performance: accuracy, confidence, and prediction of accuracy. Although evaluated using different means, *accuracy* is almost always measured as an ordinal value indicating how correct a query is. Prior studies have evaluated it using subjective assessments (Borthick et al. 2001; Chan et al. 1993) and objective assessments (e.g., Kim and March 1995; Sinha and Vessey 1999). We use an objective assessment, *semantic correctness*, defined as the ratio of required semantic elements included in a subject's query to the total number of semantic elements required in a correctly formulated query (see Appendix A). As in prior studies, the second dependent variable, *confidence*, is self-reported and measured using a five-point Likert scale.

The third dependent variable, *prediction of accuracy*, has not been examined extensively in prior studies on query formulation. It deals with how well users' self-reported confidence in the accuracy of a query predicts that query's accuracy. Without distinguishing between overconfidence and under confidence, it reports a user's absolute proficiency at self-assessment. Such a measure is important because users who exhibit higher proficiency at self-assessment more accurately

identify when queries produce the intended results and when they do not. Accordingly, they are more likely to *appropriately* rely on query results in decision-making and other managerial tasks. Although the effects of overconfidence may differ from the effects of under confidence, either can be detrimental to the effective use of organizational information resources.

Prediction of accuracy is measured using mean prediction score, a simple modification of the Mean Probability Score (Yates 1990) to account for prediction of a continuous rather than dichotomous variable. We developed this measure to compensate for problems using correlation-based measures of the proficiency of self-assessment. Specifically, when subjects express the same confidence for each prediction, a correlation between confidence and accuracy cannot be calculated. Mean prediction score does not suffer from this limitation. Mean prediction score is bounded by zero and one, with zero indicating perfect prediction; its calculation is described in Appendix B.

## Research Hypotheses

We rely on several points from our prior discussion of human cognitive processing in the formulation of this study's hypotheses.

1. Humans have a specific mental capacity for processing and recalling events in addition to the capacity for processing and recalling facts.

2. One of the primary ways that humans make sense of that with which they are not familiar is through event/narrative sense-making.

3. Temporally sequenced events are an effective mode of communication, in written forms as well in spoken forms.

Conjecturing that sense-making competencies are evoked as individuals interact with a database, we hypothesize that users will perform better at query formulation when using conceptual and logical data models that expressly represent events than they will when using conceptual and logical data models that focus on things and their states. We recognize that other cognitive processes are involved in understanding database schemata; however, other things equal, we expect that the direct representation of events will lead to a better understanding of the database, which will be manifest in better query performance (accuracy, confidence, and prediction). Accordingly, we state our hypotheses as follows:

**H1 (accuracy)**: Individuals using event-based models will formulate queries that are more accurate (semantically correct) than will individuals using state-based models.

**H2 (confidence)**: Individuals using event-based models will express higher confidence than will individuals using state-based models.

**H3 (prediction)**: Individuals using event-based models will express confidence that better predicts the accuracy of their queries than will subjects using state-based models.

Another common measure used in query performance studies is the time subjects took to compose individual queries (Chan et al. 1993). Typically, experiments are time restricted so the variance on any time-based measures is constrained. However, in our experiment, subjects were allowed to spend as much or as little time as desired. The experiment was given as an extra-credit homework assignment, and because our experience has indicated widely varying time in the completion of homework, we expected the variation in time to be too large to produce significant results. Accordingly, we do not formally hypothesize about time.

### Independent Variable and Covariates

As illustrated in Figure 2 the independent variable is the ontological foundation of the conceptual modeling method used to produce the data model, state-based, event-based, or artifact-based. Because subjects must understand the semantics of the database as well as its logical structure before they can successfully formulate SQL queries, each treatment includes both abstraction levels. The conceptual level is expressed using E-R diagrams. The logical level is expressed

using user views in the relational model. That is, an E-R diagram and its corresponding user view operationalize each treatment of the independent variable.

A single real-world domain is used in this study: the sales/collection business process of a company named TechSupport (see Appendix C). The processes and data obtained from TechSupport's operations are real, not contrived for the purpose of experimental evaluation. Thus, the experimental setting yields a high degree of realism with respect to the data subjects are asked to query. The E-R diagrams (conceptual-modeling scripts) for each treatment express similar semantics surrounding this business domain. The corresponding user views (logical-modeling scripts) are built on a common relational database copied directly from TechSupport's operational system (see Appendix D).

The distinction among treatments is rooted in the ontological status given to events. As discussed above, a common ontological position states that events and things are uniform in that they both have identity and properties (Brody 1980; Feibelman 1951; Sowa 1999; Tiles 1981). Another common ontological position states that events are changes in the states of things and, as such, cannot have properties (Bunge 1977; Wand and Weber 1995).

When the ontology underlying the conceptual modeling method views things and events uniformly, both are appropriately represented using the same construct. In an E-R conceptual-modeling grammar, the *entity* construct is appropriately used to represent events as well as things. If, however, the ontology underlying the conceptual modeling method prevents events from having properties, the *entity* construct cannot be used to represent events.

Based on this fundamental difference in ontological perspective, three treatments were developed: state-based, event-based, and artifact-based (Figures 3 through 5, respectively, and Appendix D). The *state-based* treatment conforms to an ontological position that denies properties to events. Accordingly, each of the entities in its E-R diagram represents a thing. Events are represented by changes in the values of attributes and relationships. Its development is consistent with traditional data modeling methods presented in information systems curricula (e.g., Hoffer et al. 2002) and research (e.g., Rosenthal and Reiner 1994; Teorey et al. 1986). It departs from strict adherence to the principles of Bunge's ontology because it has optional properties, which can result in null values in the database. However, the current study includes only completed sales orders, meaning that although null values are allowed, the database contains none.
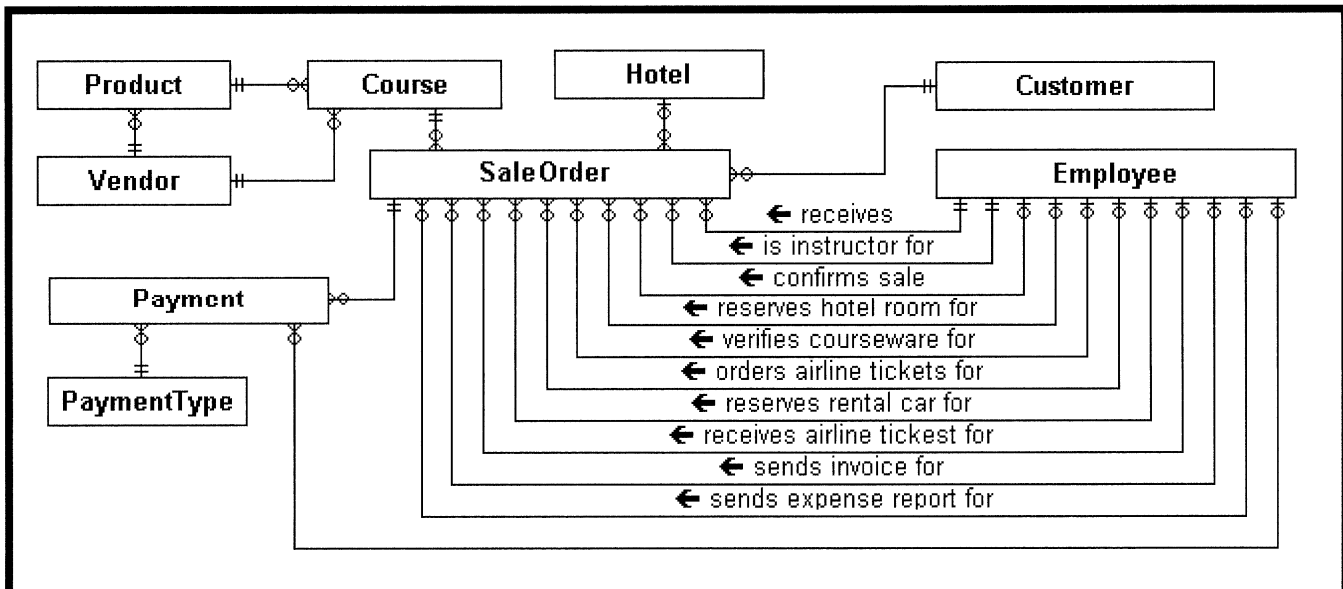
**Figure 3. State-Based E-R Diagram**

This E-R diagram is paired with a *state-based user view* of TechSupport's relational database that corresponds directly to it (Appendix D). That is, for each entity in the diagram, there is a relation in the user view with the same name. Each relationship in the diagram has a corresponding primary key/ foreign key pair in that user view (i.e., a *foreign key* column in one table that references the *primary key* column in the related table). Conversely, all primary key/foreign key pairs in that user view have corresponding relationships in the diagram.

This treatment is important because its E-R diagram is a direct mapping from the actual schema of the TechSupport database. It is the reverse-engineering equivalent to the mapping rules prescribed for constructing a normalized relational database from an E-R diagram having only binary relationships (Chiang et al. 1994).

The second treatment is *event-based*. Its E-R diagram is shown in Figure 4. It conforms to an ontological position that treats things and events uniformly. Hence, things and events are both represented using the entity construct. Its development is consistent with data modeling methods presented in accounting information systems curricula[3] (Denna et al. 1993; Hollander et al. 1999) and research (Geerts and McCarthy 2002; McCarthy 1982). As with the state-based treatment, the

event-based E-R diagram is paired with an *event-based user view* of TechSupport's database. It has a similar one-to-one mapping between entities and relations and between relationships and primary key/foreign key pairs (Appendix D).

The event-based treatment is important because it is the extreme representation of events. The sales/collection process at TechSupport has one event that receives an order, one event that confirms the sale, one that reserves a hotel, one that orders flight tickets, one that reserves a car, one that verifies courseware, one that receives flight tickets, one that sends an invoice, and one that sends an expense report (see Appendix C). Each of these events is represented as an entity. ReceiveOrder corresponds to the entity SaleOrder in the state-based representation. Both are identified by the attribute SaleID. However, while SaleOrder contains all attributes that are functionally dependent on SaleID, ReceiveOrder only contains attributes that are relevant to that event (see Appendix D).

The additional entities in the event-based model are related one-to-one to ReceiveOrder; hence they each share its identifier, SaleID. Each has attributes and relationships that are relevant to its respective event. VerifyCourseware, for example, has the attributes, CoursewareHandled (how the courseware was handled—"sent to TechSupport" or "sent to client") and CoursewareVerified (the date the courseware was verified to have arrived) and a relationship with the employee entity, indicating who verified that the courseware was ready for instruction.
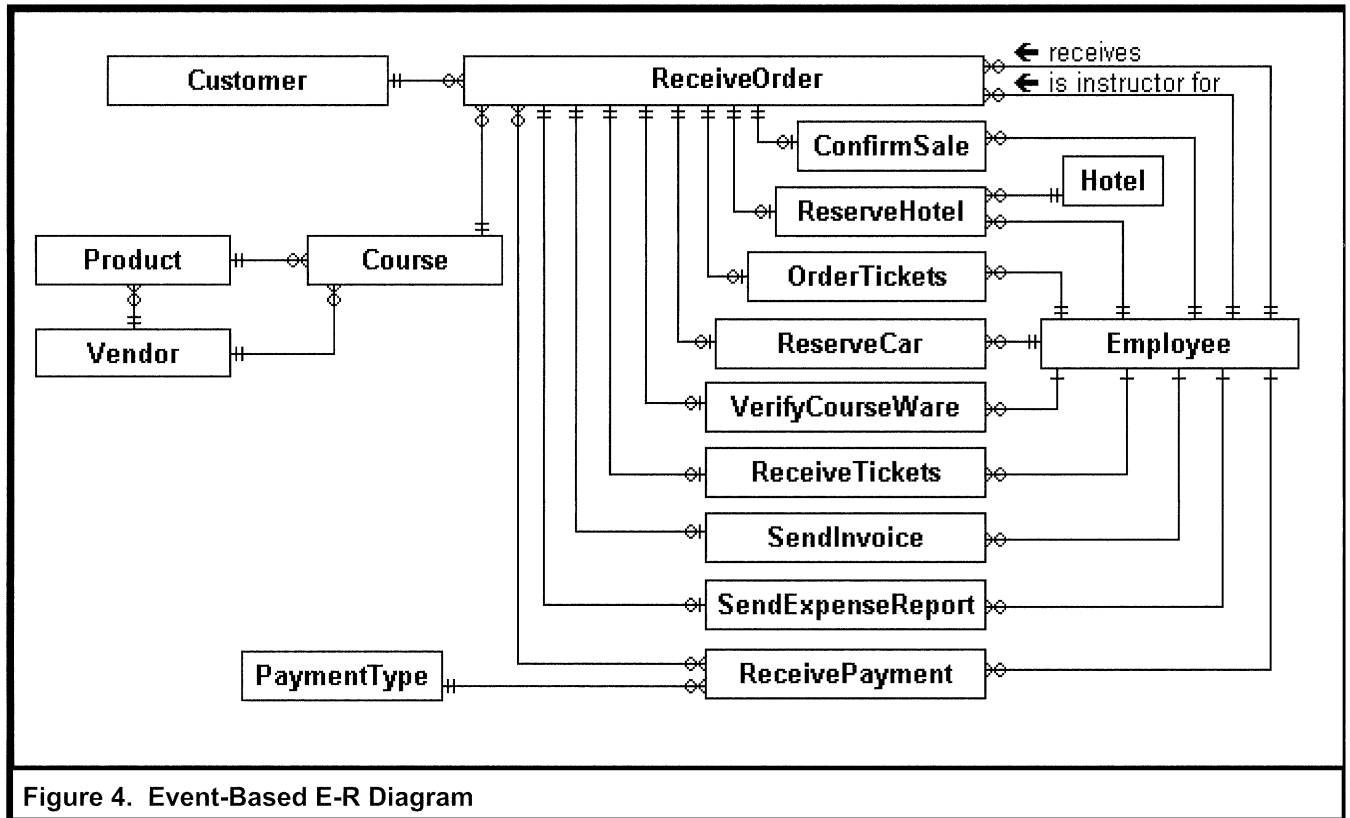
---

[3]This diagram does not strictly adhere to the conventions of REA, which specify the representation of the "give-take" duality that characterizes relationships among economic events.

**Figure 4. Event-Based E-R Diagram**

We note that the influence of such diverse ontological positions in a requirements elicitation task would likely lead to the expression of different domain semantics. However, the development of the treatments was constrained to express the same semantics recorded in TechSupport's existing database. This constraint is important because this study examines the effects of the independent variable on query formulation, not on requirements elicitation. Accordingly, the treatments must differ in the *way* in which domain semantics are conveyed without being confounded by the representation of *different* domain semantics. Hence while the event-based treatment (Figure 4) has an *entity* for each event, an instance of which is created when the corresponding business activity is completed, the state-based treatment (Figure 3) has a *relationship* for each event, an instance of which is created when the corresponding business activity is completed.

The state-based treatment organizes information around things, while the event-based treatment organizes information around events. This difference is seen both in the clustering of attributes and in the naming conventions of the entities. A positive finding in the comparison of these two treatments would leave an important question unanswered. Is the treatment effect a result of clustering attributes around events

or is it a result of linguistic choices in the naming of entities ("Payment" versus "ReceivePayment")? To help answer this question, we developed a third treatment, which we term *artifact-based*. It is state-based in the sense that the entities in the diagram are named after things; however, it is event-based in the sense that the things around which attributes are clustered are organizational forms and documents (artifacts) used to record data about events.[4] The performance of subjects presented with this treatment will help us understand if any treatment effect between the event-based treatment and state-based treatment is linguistic or structural or a combination of both.

The E-R diagram for the artifact-based treatment is shown in Figure 5. It has a corresponding *artifact-based user view* (Appendix D). It captures the common practice of constructing entities for named artifacts (forms and documents). Organizations often create such artifacts to give visibility and prominence to abstractions such as events that are important elements of business processes.

---

[4]Note that these would not be considered "things" in Bunge's ontology (Wand and Weber 2002).
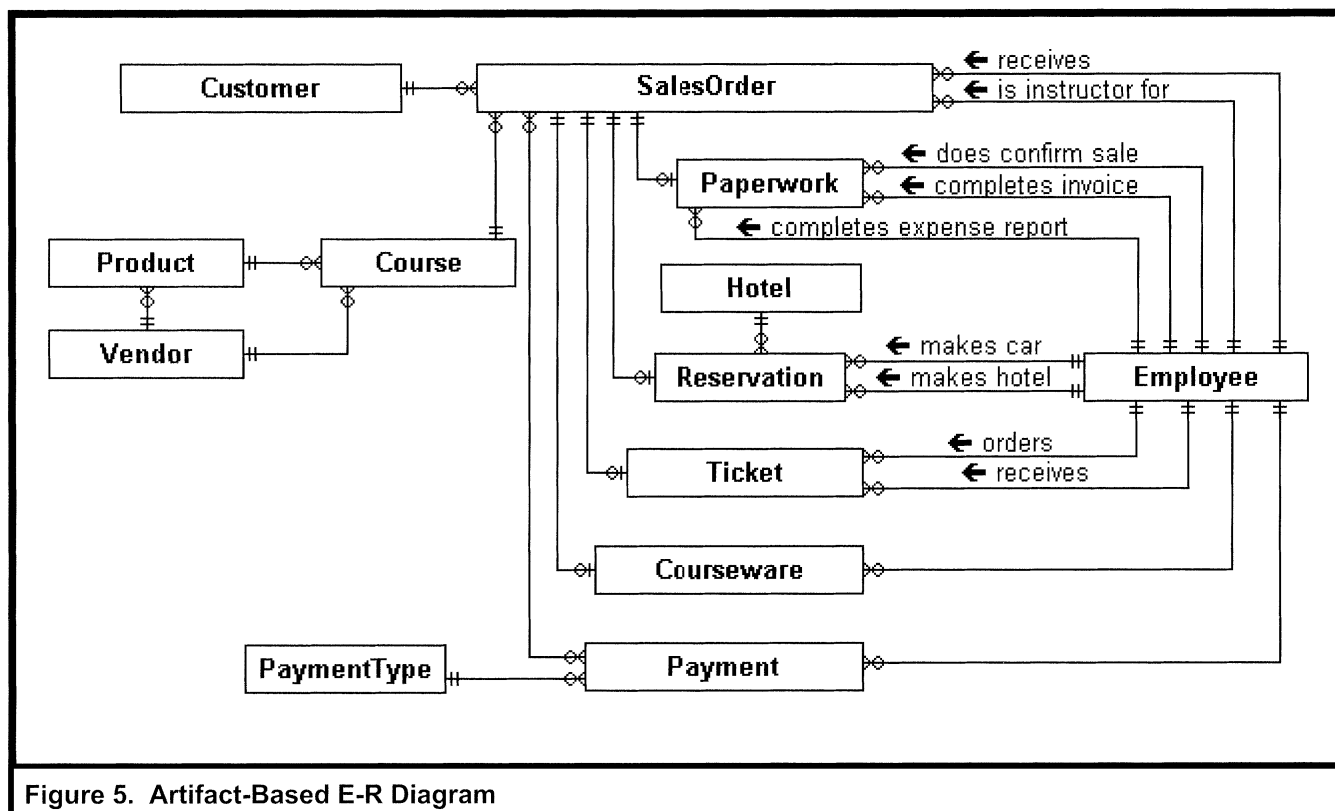
**Figure 5. Artifact-Based E-R Diagram**

While this treatment is constructed to record data about events, it differs from the event-based treatment in several important ways. First, the naming convention is document-focused. Entities are named for the documents used in the organization and do not specifically convey events that can be composed into a narrative. Second, not all events have corresponding documents and organizations may record several events on the same document. In fact, this treatment has only one entity that directly corresponds to an event, Courseware, corresponding to the event Verify Courseware. It has three entities that correspond to combined events. Paperwork corresponds to the combination of Confirm Sale, Send Invoice, and Send Expense Report; Reservation corresponds to the combination of Reserve Hotel and Reserve Car; Ticket corresponds to the combination of Order Tickets and Receive Tickets. It forms a middle ground between the state-based and event-based treatments with respect to complexity and to the representation of events.

## Research Procedures

A one-factor between-subjects experiment was conducted to investigate the effects of the independent variable (ontological foundation) on the dependent variable (query performance). Because the experiment requires subjects to read and interpret an E-R diagram, the potential for the researchers to create a training bias in the subjects is a significant threat to validity. Accordingly, the researchers were not involved with the training of subjects. Moreover, it was decided that the subject pool should be drawn from several different educational backgrounds to reduce the likelihood that subjects' pre-experiment training unduly favored any treatment.

To meet these constraints, an instrument was constructed to conduct the experiment via the Internet. The instrument was built using a combination of HTML, Active Server Pages, and JavaScript. A pilot study involving 36 subjects was conducted to test the instrument and refinements were made prior to conducting the experiment.

Subjects were recruited from information systems programs at six universities in North America and Europe and randomly assigned to treatments. Subjects were either enrolled in an introductory database management course, or in a course for which database management was a prerequisite. Course extra credit was given for participation and 342 subjects produced useable results. In an exit survey, subjects disclosed their

approximate age, ethnic background, gender, academic major, comfort level in reading E-R diagrams, and comfort level in writing SQL queries. An analysis of the subject assignment showed no systematic bias in any treatment group on any of these measures. *Comfort in reading E-R diagrams* and *comfort in writing SQL queries* were considered as covariates. Because they held substantial correlation (coefficient of correlation 0.56) comfort in reading E-R diagrams was selected to serve as the model's covariate.

All subjects used the same instrument. It provides a textual description of the TechSupport business activity (Appendix C), an E-R diagram and user view corresponding to one treatment, a place to formulate SQL queries and execute them against their user view, a place to see the results of their queries when executed (or error messages for syntactically incorrect queries), a help system that includes information about the conceptual-modeling grammar used in the study as well as help on SQL syntax, and a set of information requests that define the experimental task (Appendix E).

The information requests were developed so as not to favor one treatment over another. Although it is possible to write an information request that requires subjects in one treatment to formulate a more complex query than subjects in another treatment (e.g., by requiring additional joins), all were composed to ensure a similar level of difficulty for all treatments. They were presented to three experts in database management, each an author of a different collegiate database management textbook. These authors were asked if subjects in one treatment would have an advantage over subjects in other treatments in building queries to answer the information requests. They saw no treatment advantage for any information request. Examples of correct queries for each treatment are shown in Appendix E.

Subjects were given the URL for the study and an access code and asked to formulate an SQL query for each of the information requests at a time and place of their choosing. They could stop the task and start again where they left off at their discretion. They could view the E-R diagram appropriate for their treatment, go back to the textual description, execute trial queries, and view help for SQL syntax. When subjects clicked on one of the entities in their E-R diagram, they were presented with the logical description of the relation corresponding to that entity, including a description of the data held in each of the attributes (see Appendix D). When subjects were satisfied that their SQL query appropriately fulfilled an information request, they indicated their confidence in the accuracy of their answer on a one to five scale. All queries and corresponding confidence assessments were stored for automated analysis. Subjects could freely

refer to different parts of the instrument and could change both the queries they submitted and their confidence in the accuracy of those queries.

The realism of enabling subjects to explore the database, execute their queries, see the results and iteratively revise and re-execute their queries is important. Prior studies that focus on constructing SQL queries without this iterative ability miss a significant component of the human cognitive processes required to effectively utilize data resources. This is well evidenced in our log files. Subjects commonly examined the data in individual tables while constructing multi-table queries and frequently revised a query based on the results obtained from an initial attempt.

# Results

An analysis of variance (ANOVA) test of significance was conducted (Table 1). In addition to using the treatment as a predictor, subjects' self-reported comfort level in reading E-R diagrams was used as a model covariate. This analysis did not indicate a significant treatment effect on query accuracy (H1). In fact, the mean percentage of semantic elements identified (semantic correctness) ranges only from 87 to 88 percent across all treatments. Further, the results showed no treatment effect on the confidence users expressed regarding the accuracy of their queries (H2). The average confidence ranged only from 3.90 to 4.02 on a five-point scale. However, the results indicated that the treatment did have a significant effect on prediction of accuracy (H3).

Noting that 0 is a perfect score for prediction of accuracy, the event-based treatment exhibited the best score (.058); the artifact-based treatment exhibited the worst score (.071); and the state-based treatment exhibited a score between them (.065). The p-values for the pairwise mean comparisons among treatments indicate that there is a significant difference in subjects' prediction of accuracy between the state-based treatment and the event-based treatment ($p = 0.0210$) and between the event-based treatment and the artifact-based treatment ($p = 0.0041$). However, they do not indicate a significant difference between the state-based treatment and the artifact-based treatment ($p = 0.6394$).

For this finding, the interaction term between the predictor variable and the covariate, comfort level with reading E-R diagrams, is also significant (0.0176). This indicates that the effect of the treatment is different for various levels of the covariate. Figure 6 shows the interaction chart for the three treatments and various levels of comfort in reading E-R dia-

**Table 1. Experimental Results**

| Hypotheses | Treatment Mean[†] | | | ANOVA P-Values | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Pairwise Mean Comparisons | | |
| | S | E | A | Treatment | Covariate | Interaction | S-E | S-A | E-A |
| H1: Accuracy[‡] | 0.88 | 0.88 | 0.87 | 0.2592 | < 0.0001** | 0.1775 | | | |
| H2: Confidence[††] | 4.02 | 4.01 | 3.90 | 0.1296 | < 0.0001** | 0.1994 | | | |
| H3: Prediction[‡‡] | 0.065 | 0.058 | 0.071 | 0.0076 | < 0.0001** | 0.0176* | 0.0210* | 0.6394 | 0.0041** |

*significant at alpha = 0.05    **significant at alpha = 0.01
[†]S = State Based, E = Event Based, A = Artifact Based
[‡]Average across queries, maximum is 1.
[††]Confidence is measured as a five-point, self-reported Likert variable, averaged across queries.
[‡‡]Prediction of accuracy is measured using mean prediction score, 0 equals perfect prediction.
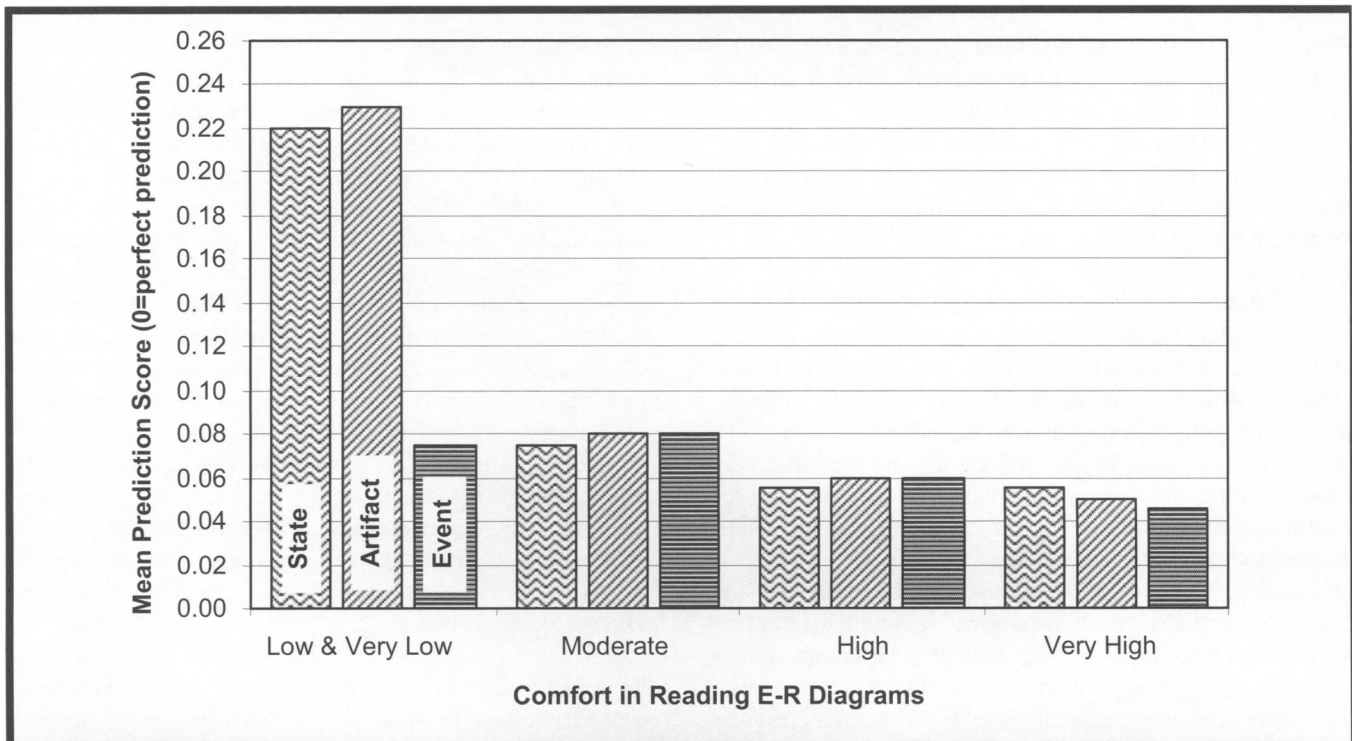


**Figure 6. Interaction Between *Comfort in Reading E-R Diagrams* and *Treatment* on Prediction of Accuracy**

grams on subjects' prediction of accuracy. In general, as subjects' comfort in reading E-R diagrams increases, so does their prediction of accuracy (i.e., the mean prediction score decreases).

For subjects expressing low or very low comfort in reading E-R diagrams, their prediction of accuracy was worse (i.e.,

the mean prediction score was higher) than for subjects expressing moderate, high, or very high comfort in reading E-R diagrams for the state-based and the artifact-based treatments. However, for the event-based treatment, the prediction of accuracy was about the same for all subjects independent of expressed comfort in reading E-R diagrams.

| Table 2. Analysis of Treatment Effects on Time Spent | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | **Treatment Mean**[†] | | | **ANOVA P-Values** | | | | |
| | | | | | | **Pairwise Mean Comparisons** | | |
| ***Post Hoc* Analysis** | **S** | **E** | **A** | **Treatment** | **Covariate** | **S-E** | **S-A** | **E-A** |
| Total time spent on experiment | 86.2 | 81.1 | 72.3 | 0.0185* | 0.6314 | 0.3850 | 0.0061** | 0.0548 |
| Time spent on semantics[‡] | 54.0 | 51.8 | 46.3 | 0.0312* | 0.8021 | 0.7429 | 0.0166** | 0.0343* |
| Time spent on queries 3-7 | 14.2 | 10.6 | 9.8 | < 0.0001** | 0.9073 | 0.0004** | < 0.0001** | 0.4329 |

*significant at alpha = 0.05          **significant at alpha = 0.01
[†]Measured in minutes; S = State Based, E = Event Based, A = Artifact Based
[‡]Time subjects spent viewing data representation, viewing the field structure of underlying tables, and reading the textual description of the business process.

Because no significant treatment effect was found for hypothesis 2 (confidence), we conducted a *post hoc* analysis of its corollary: did the treatment have an effect on how long subjects spent to achieve the confidence level at which they were comfortable moving on? This question was examined in a similar manner to the tests of hypotheses (Table 2) using data from system log files.

The system log files contain a time stamped entry for each action of each subject. This action-time log provides the ability to tell when and for how long each page of the instrument was displayed. However, it cannot tell how long subjects were actually engaged with the displayed page. To control for the possibility of subjects leaving the instrument for several minutes, page-view times were truncated at three minutes. This threshold was chosen after observing additional individuals complete the experimental task in a laboratory setting. Under these conditions, no subject engaged any page view for more than two consecutive minutes (returning to the same page after visiting another page is logged as a different page-view). Recognizing differences between laboratory and experimental settings, three minutes was deemed appropriate. The analysis of time spent was conducted with various thresholds between 2 and 10 minutes without a material difference in the results; 97 percent of the page-view intervals were shorter than 2 minutes. Accordingly, the means reported in Table 2 are not affected by long periods of user inactivity.

The results shown in Table 2 indicate that the treatment had a significant effect on the overall time subjects spent completing the experiment (p = 0.0185) but the covariate, comfort level in reading E-R diagrams, did not (p = 0.6314). The average time spent by subjects in the study was highest for the state-based treatment (86.2 minutes) and lowest for the artifact-based treatment (72.3 minutes). However, the variation is such that the only significant contrast is between the state-based and artifact-based treatments (p = 0.0061). We further segmented the time subjects spent into those activities directly related to gaining an understanding of semantics from all other activities in the experiment. To calculate the time spent on semantics, we summed the time a subject spent showing the E-R diagram, the time spent showing the names and descriptions of fields in the database schema, and the time spent showing the textual description of TechSupport's business process. Other activities such as looking in help, viewing error messages, viewing the results of a query, or reading information requests were likely working on understanding syntax, administrative activities (as required by the experimental instrument), or in the process of verifying the reasonableness of the query results for the information request. As such, they were considered to be activities less directly indicative of subjects engaging in sense-making processes.

Examination of time spent on semantics indicates a significant treatment effect (p = 0.0312). The comparison between treatments indicates no significant difference between the state-based and event-based treatments (p = 0.7429). It indicates a significant difference between the other pairs: the artifact-based treatment is significantly less than both state-based (46.3 minutes compared to 54.0 minutes; p = 0.0166) and event-based treatments (51.8 minutes; p = 0.0343). Finally, the amount of time spent on queries three through seven was considered. The first two queries deal exclusively with portions of the E-R diagrams and database schemata that were identical across treatments. They were positioned at the beginning of the experiment to allow subjects to become comfortable with the instrument before beginning the portion of the experiment that was directly affected by the treatments.

For these queries, the treatment effect is significant ($p <$ .0001). The p-values for the pairwise mean comparison among treatments indicate no significant difference between the event-based and artifact based treatments (10.6 minutes and 9.8 minute, respectively; $p = 0.4329$). However, each is significantly lower than the state-based treatment (14.2 minutes; $p = 0.0004$ and $p < 0.0001$, respectively).

## Discussion and Conclusion ▬▬▬▬▬

This study found no significant difference among treatments for hypothesis 1 (accuracy). This hypothesis called for the treatments to be examined in aggregate over each of the seven information tasks. A *post hoc* analysis of the accuracy of subjects' queries was conducted for each information task individually. Like the aggregate results, there were no significant differences in accuracy across treatments for any of them. The statistical power of these tests is not sufficient to conclude that a treatment effect does not exist; rather, that we were unable to observe one.

One possible explanation for this result is the complexity of the diagrams. The state-based treatment is least complex (9 entities, 19 relationships), the event-based treatment is most complex (17 entities, 27 relationships), and the artifact-based treatment is between them (13 entities, 23 relationships). It may be that the cognitive benefits of the event-based treatment were able to offset the effects of increased complexity but could not overcome them. In our study, several of the events of the business process had the characteristic of being related to each other in a one-to-one fashion, thus enabling them to be represented by a single entity in the state-based treatment. Although this may be a common characteristic of diagrams produced using an event-based conceptual-modeling method, it is not a universal one. Thus it is possible that event-based E-R diagrams of equal complexity to similar state-based E-R diagrams would lead to superior accuracy in query formulation tasks. Beyond the complexity issue, there are several other possible reasons for lack of support for hypothesis 1.

In the *ad hoc* query formulation process, our treatments should only affect the mapping of terms in the information request to elements of the database schema. Hypothesis 1 relies on the ability of the treatments to provide different levels of understanding of the underlying content of the database. Although our study references a business process with which subjects likely had little or no experience, it uses a straightforward implementation of a business' sales/collection process. It is possible (even likely) that subjects were able to

effectively apply what they knew about other sales/collection processes to the experimental task in such a way that they were easily able to make sense of the database schema independent of the treatment. Like the narrative sense-making competency, the transfer of learning from one context to another distinct, yet similar, context is thought to be a fundamental component of human cognitive processes (Shepard 1987).

Another possible explanation for not finding support for hypothesis 1 may stem from the information requests used in the study. For a treatment effect to be observed, subjects would need to engage their sense-making capabilities in the experimental task. If information requests were worded in such a way that subjects could reasonably translate them into table and column names without needing to grasp the underlying semantics, then the treatments would not be expected to have differential effects. Subjects may have performed a lexical mapping from terms in the information requests to the names of entities, attributes, and relationships (tables and columns) in the data dictionary (Appendix D) without needing to engage their sense-making capabilities. Although this possibility was considered in the formulation of the information requests, the ability of subjects to accomplish lexical mappings without domain understanding may have been underestimated. One area of potential future research is to examine the performance of writers of *ad hoc* queries under conditions of context-free information requests, familiar-context information requests, and counterintuitive information requests. This kind of study would further the understanding of balances and tradeoffs between lexical mapping and domain understanding.

We examined several measures of the processes subjects used to complete the experimental task to see if there was some other indication of significant differences among the treatments. We studied how long subjects spent in different parts of the experimental instrument, how many queries they tested as they worked to complete the task, and the number of times they switched among the different pages of the instrument. The analysis of time spent gave some interesting insights, and is discussed below. However, all other analyses of process showed no significant treatment effects.

This study found no significant difference among treatments for hypothesis 2 (confidence). We predicted higher confidence for the event-based treatment based on the reasoning that (1) although the experiment was not time constrained, competing demands on subjects' time would compel them to move on from answering an information request before they reached complete confidence in their response and (2) if the event-based treatment resulted in a deeper understanding of

the database schema, then subjects in this treatment would reach higher confidence before time constraints compelled them to move on.

When we found no significant difference in confidence expressed, we examined how long subjects took to achieve that level of confidence. The results of the *post hoc* analysis of time spent (Table 2) indicate that considering (1) time spent on the experiment, (2) time spent on semantics, and (3) time spent on queries three through seven that the state-based treatment is always dominated by a treatment that gives prominence to events, either by their direct representation or by the representation of their associated artifacts. For two of the three time elements examined, the artifact-based treatment produced significantly better performance than did the event-based treatment. We conjecture that this result is due to the increased complexity of the event-based treatment compared to the artifact-based treatment. However, when isolating the comparison to just those queries that were directly affected by the treatment, the event-based and the artifact-based treatments both lead to superior performance as compared to the state-based treatment, and did not lead to significantly different performance as compared to each other.

The analysis of time spent on queries three through seven is interesting for three reasons. First, it shows that the vast majority of the total time spent (84 to 88 percent) was dedicated to subjects gaining an understanding of the instrument, task, and data models well enough to complete the first two queries. Second, once that initial level of understanding was attained, the treatment effect clearly shows a reduction in time spent formulating queries when events are represented, either directly or through their associated artifacts. This suggests that the additional complexity exhibited by the event-based treatment over the state-based treatment negatively affects time spent initially understanding the models but that the added complexity matters less as users become familiar with them. Furthermore, even though the state-based treatment exhibited the least complexity, its subjects took more time to complete this segment of the study, indicating that once users gain some familiarity with the diagrams, the treatment effect is strong enough to overcome the complexity effect. This finding is not dependent on subjects' comfort level in reading E-R diagrams, but holds true for all comfort levels. Third, this finding suggests that the linguistic effect of naming entities to represent business events is of little consequence. Whether an entity was named as a noun (i.e., *payment* or *courseware*) or it is named as a verb (i.e., *receive payment* or *verify courseware*) seems to be not nearly as important in the sense-making process as the fact that information pertaining to business events is given prominence in the model.

A related factor that may have contributed to the lack of support for hypothesis 2 is found in the manner in which the experiment was administered. Since subjects from six different universities in North America and Europe participated in the experiment, setting up similar laboratory conditions for all subjects was improbable. Instead, subjects completed the experiment through the Web interface at a time and place of their convenience. Although this reduced the control of the experimental setting, it was deemed an acceptable tradeoff against the potential of introducing a systematic bias through differing laboratory settings. The result is that subjects were able to trade time spent on the experiment for other activities. This gives some insight into the process of *ad hoc* query formulation. Although the above analysis leads us to believe that the hypothesis of higher confidence in event-based treatments would be supported in a time-constrained study, this finding would have little meaning for the business environment because such time constraints are artificial. It seems that query writers will take the time needed to reach a confidence level with which they are comfortable, but that they reach that comfort level more quickly working with representations in which events are expressly represented.

Hypothesis 3 (prediction) was supported for the class of subjects who reported low or very low levels of comfort in reading E-R diagrams. Such subjects in the event-based treatment expressed confidence that better predicted the accuracy of their queries than did such subjects in the state-based or artifact-based treatments. This finding is made all the more compelling by the lack of support for hypotheses 1 and 2. If significant treatment effects were found either for accuracy or for confidence, there would be some question as to whether support for hypothesis 3 is due to the treatment or is simply a result of higher confidence or improved accuracy. However, when subjects demonstrate similar accuracy and confidence and still demonstrate superior prediction of accuracy, there can be no question that the treatment holds a direct effect.

This finding is consistent with prior work examining the connection between human cognition and the use of conceptual database schema. Ramesh and Browne (1999) found that causation was better expressed by individuals with little or no exposure to data modeling principles than by individuals with prior exposure. Although their study examined deficiencies in the E-R model for representing causality, their finding is of particular interest because, like the current study, it shows that human information processing competencies surrounding causality have particular importance for individuals with little database experience.

Hence we conclude that an event-based E-R diagram can lead casual users to more accurately recognize when queries they

have formulated are correct. Failure to accurately recognize when queries are accurately formulated could lead to significant judgmental errors and improper decisions. Because we evaluate this effect using a new measure (mean prediction score), more research is needed to understand how its assessment relates to user behavior in real-world settings.

We further conclude that the use of E-R diagrams and corresponding user views that give prominence to some representation of events allow users to more quickly formulate queries without sacrificing accuracy or confidence. These findings are particularly compelling in a business environment where managers seek to make sense of transaction (event) data through data mining and business intelligence interfaces. In such environments, E-R diagrams are frequently used to communicate the structure of organizational data to end-users (MicroStrategy 2003).

It should be pointed out that none of the subjects expressly received training in event-based conceptual-modeling methods. A question that arises is, can we develop strong methods for reading graphical data models that leverage the human competency for processing events and narratives? This is an area that merits further research. Moreover, this experiment was conducted using event, state, and artifact-based data models that are nearly logically equivalent. Although it is possible to model temporal semantics in such a way that an event-based data model and a state-based data model will convey nearly identical information, it is not likely that the natural application of these two ontologically diverse methods will yield data models that equivalently represent a given domain. That is, when an analyst models a domain using a conceptual modeling method that treats things and events uniformly, allowing events to have attributes, he or she will likely record different semantics than will an analyst using a method that treats things and events differentially, not allowing events to have attributes.

The findings of this research have demonstrated that the human competency for processing events can be leveraged though data models that emphasize events—even in the constrained context of querying a common underlying relational database. Future research will study the effects of event-based and state-based data models (grammars, scripts, and methods) on human performance in more conceptual tasks such as information system requirements determination and validation.

## Acknowledgments

## References

Antony, S. R., and Batra, D. "CODASYS: A Consulting Tool for Novice Database Designers," *ACM SIGMIS Database* (33:3), Summer 2002, pp. 54-68.

Batra, D., Hoffer, J. A., and Bostrom, R. P. "A Comparison of User Performance Between the Relational and the Extended Entity Relationship Models in the Discovery Phase of Database Design," *Communications of the ACM* (33:2), February 1990, pp. 126-139.

Bodart, F., Patel, A., Sim, M., and Weber, R. "Should the Optional Property Construct Be Used in Conceptual Modeling? A Theory and Three Empirical Tests," *Information Systems Research* (12:4), 2001, pp. 384-405.

Borthick, A. F., Bowen, P. L., Jones, D. R., and Tse, M. H. K. "The Effects of Information Request Ambiguity and Construct Incongruence on Query Development," *Decision Support Systems* (32), 2001, pp. 35-56.

Bowen, P. L., O'Farrell, R. A., and Rohde, F. H. "How Does Your Model Grow? An Empirical Investigation of the Effects of Ontological Clarity and Application Domain Size on Query Performance," in *Proceedings of the 25th International Conference on Information Systems*, R. Agarwal, L. Kirsch, and J. I. DeGross (eds.), Washington, DC, December 12-15, 2004, pp. 77-90.

Brody, B.A. *Identity and Essence*, Princeton University Press, Princeton, NJ, 1980.

Bronts, G. H. W. M., Brouwer, S. J., Martens, C. L. J., and Proper, H. A. "A Unifying Object Role Modeling Theory," *Information Systems* (20:3), 1995, pp. 213-235.

Bunge, M. *Ontology I: the Furniture of the World*, Volume 3 of *Treatise on Basic Philosophy*, D. Reidel Publishing Company, Boston, MA, 1977.

Burton-Jones, A., and Weber, R. "Understanding Relationships with Attributes in Entity-Relationship Diagrams," in *Proceedings of the 21st International Conference on Information Systems*, P. De and J. I. DeGross (eds.), Charlotte, NC, December 13-15, 1999, pp. 214-228.

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

Chen, P. P. S. "The Entity-Relationship Model-Toward a Unified View," *ACM Transactions on Database Systems* (1:1), 1976, pp. 9-36.

Chiang, R. H. L., Barron, T. M., and Storey, V. C. "Reverse Engineering of Relational Databases: Extraction of an EER Model from a Relational Database," *Data and Knowledge Engineering* (12), 1994, pp. 107-142.

Denna, E. L., Cherrington, J. O., Andros, D. P., and Hollander, A. S. *Event-Driven Business Solutions: Today's Revolution in Business and Information Technology*, Business One Irwin, Homewood IL, 1993.

Dey, D., Barron, T., and Storey, V. "A Conceptual Model for the Logical Design of Temporal Databases," *Decision Support Systems* (15), 1995, pp. 305-321.

Feibleman, J. K. *Ontology*, The Johns Hopkins Press, Baltimore, MD, 1951.

Gee, J. P. "The Narrativization of Experience in the Oral Style," *Journal of Education* (167), 1985, pp. 9-35.

Geerts, G. L., and McCarthy, W. E. "An Ontological Analysis of the Economic Primitives of the Extended-REA Enterprise Information Architecture," *International Journal of Accounting Information Systems* (3:1), March 2002, pp. 1-16.

Gemino, A., and Wand, Y. "Complexity and Clarity in Conceptual Modeling: Comparison of Mandatory and Optional Properties," *Data and Knowledge Engineering* (55:3), December 2005, pp. 301-326.

Goodhue, D. L., Klein, B. D., and March, S. T. "User Evaluations of IS as Surrogates for Objective Performance," *Information and Management* (38), 2000, pp. 87-101.

Gray, P., and Watson, H. J. *Decision Support in the Data Warehouse*, Prentice Hall, Upper Saddle River, NJ, 1998.

Halevy, A. Y. "Answering Queries Using Views: A Survey," *The VLDB Journal* (10), 2001, pp. 270-294.

Halpin, T. A. *Information Modeling and Relational Databases: From Conceptual Analysis to Logical Design* (3rd ed.), Morgan-Kauffmann, San Francisco, 2001.

Hoffer, J. A., Prescott M. B., and McFadden, F. R. *Modern Database Management* (6th ed.), Prentice Hall, Upper Saddle River, NJ, 2002.

Hollander, A., Cherrington, J. O., and Denna, E. L. *Accounting, Information Technology, and Business Solutions* (2nd ed.), McGraw-Hill/Irwin, Boston, MA, 1999.

Jih, W. J., Bradbard, D. A., Snyder, C. A., Thompson, N. G. "The Effects of Relational and Entity-Relationship Data Models on Query Performance of End Users," *International Journal of Man-Machine Studies* (31:3), 1989, pp. 257-267.

Kent, W. *Data and Reality*, Elsevier Science Limited, Amsterdam, 1978.

Kim, Y. G., and March, S. T. "Comparing Data Modeling Formalisms," *Communications of the ACM* (38:6), June 1995, pp. 103-115.

Lakoff, G. *Women, Fire and Dangerous Things: What Categories Reveal about the Mind*, University of Chicago Press, Chicago, IL, 1987.

Leitheiser, R. L., and March, S. T. "The Influence of Database Structure Representation on Database System Learning and Use," *Journal of Management Information Systems* (12:4), Spring, 1996, pp 187-213.

Markowitz, V. M., and Shoshani, A. "Representing Extended Entity-Relationship Structures in Relational Databases: A Modular Approach," *ACM Transactions on Database Systems* (17:3), September 1992, pp. 423-464.

McCarthy, W. E. "The REA Accounting Model: A General Framework for Accounting Systems in a Shared Data Environment," *The Accounting Review* (57:3), July 1982, pp 554-578.

MicroStrategy, Inc. *MicroStrategy 7i Sales Analysis Module Reference Guide* (Version 7.2.3), 2003 (available online at http://www.teradatastudentnetwork.com/, accessed April 1, 2004).

Nyberg, L. "Mapping Episodic Memory," *Behavioral Brain Research* (90), 1998, pp. 107-114.

Orr, J. "Sharing Knowledge, Celebrating Identity: Community Memory in a Service Culture," in *Collective Remembering*, D. Middleton and D. Edwards (eds.), Sage Publications, Newbury Park, CA, 1990.

Parsons, J., and Wand, Y. "Emancipating Instances from the Tyranny of Classes in Information Modeling," *ACM Transactions on Database Systems* (25:2), June 2000, pp. 228-268.

Pillemer, D. B. *Momentous Events, Vivid Memories,* Harvard University Press, Cambridge, MA, 1998.

Ramesh, V., and Browne, G. "Expressing Casual Relationships in Conceptual Database Schema," *The Journal of Systems and Software* (45), 1999, pp. 225-232.

Robinson, J. A., and Hawpe, L. "Narrative Thinking as a Heuristic Process," in *Narrative Psychology: The Storied Nature of Human Conduct*, T. R. Sarbin (ed.), Prager Publishers, New York, 1986.

Rosenthal, A., and Reiner, D. "Tools and Transformations—Rigorous and Otherwise—for Practical Database Design," *ACM Transactions on Database Systems* (19:2), June 1994, pp. 167-211.

Shasha, D. "Tuning Databases for High Performance," *ACM Computing Surveys* (28:1), March 1996, pp. 113-115.

Shepard, R. N. "Evolution of a Mesh between Principles of the Mind and Regularities of the World," in *The Latest on the Best*, J. Dupre (ed.), The MIT Press, Cambridge, MA, 1987.

Siau, K. L., Chan, H. C., and Wei, K. K. "Effect of Query Complexity and Learning on Novice User Query Performance with Conceptual and Logical Database Interface," *IEEE Transactions on Systems, Man, and Cybernetics: Part A* (34:2), March 2004, pp. 276-282.

Silberschatz, A., and Korth, H. F. "Data Models," *ACM Computing Surveys* (28:1), March 1996, pp. 105-108.

Sinha, A. P., and Vessey, I. "An Empirical Investigation of Entity-Based and Object-Oriented Data Modeling: A Development Life Cycle Approach," in *Proceedings of the 20ᵗʰ International Conference on Information Systems*, P. De and J. I. DeGross (eds.), Charlotte, NC, December 12-15,1999, pp. 229-244.

Sowa, J. F. *Knowledge Representation: Logical, Philosophical, and Computational Foundations*, Brooks/Cole Publishing, Pacific Grove, CA, 1999.

Speier, C., and Morris, M. G. "The Influence of Query Interface Design on Decision-Making Performance," *MIS Quarterly* (27:3), September 2003, pp. 397-423.

Teorey, T. J., Yang, D., and Fry, J. P. "A Logical Design Methodology for Relational Databases Using the Extended Entity-Relationship Model," *ACM Computing Surveys* (18:2), June 1986, pp. 197-222.

Tiles, J. E. *Things that Happen*, Aberdeen University Press, Aberdeen, Scotland, 1981.

Tulving, E. *Elements of Episodic Memory*, Oxford University Press, New York, 1983.

Tulving, E. "Episodic Memory: From Mind to Brain," *Annual Review of Psychology* (53), 2002, pp. 1-25.

Wand, Y., Storey, V. C., and Weber, R. "An Ontological Analysis of the Relationship Construct in Conceptual Modeling," *ACM Transactions on Database Systems* (24:4), December 1999, pp. 494-528.

Wand, Y., and Weber, R. "On the Deep Structure of Information Systems," *Information Systems Journal* (5), 1995, pp. 203-233.

Wand, Y., and Weber, R. "Research Commentary: Information Systems and Conceptual Modeling: A Research Agenda," *Information Systems Research* (13:4) December 2002, pp. 363-376.

Weber, R., and Zhang, Y. "An Ontological Evaluation of NIAM's Grammar for Conceptual Schema Diagrams," in *Proceedings of the 11ᵗʰ International Conference on Information Systems*, J. I. DeGross, I. Benbasat, G. DeSanctis, and C. M. Beath (eds.), New York, December 16-18, 1991, pp 75-82.

Yates, J. F. *Judgment and Decision Making*, Prentice-Hall, Englewood Cliffs, NJ, 1990.

## About the Authors

**Gove N. Allen** holds bachelor's and master's degrees in Accountancy from Brigham Young University. Dr. Allen received his Ph.D. from the University of Minnesota in 2001 and currently serves as an assistant professor of e-business and information systems at Tulane University's A. B. Freeman School of Business. He has consulted on the implementation of database technology for many major corporations including Sony, AT&T, Sprint, Hewlett Packard, Micron, Intel, 3M, American Express, and NASA. More recently, he developed WebSQL.org, a site for teaching database management that allows dynamic execution of Structured Query Language against Oracle databases through a simple web interface.

**Salvatore T. March** is the David K. Wilson Professor of Management at the Owen Graduate School of Management, Vanderbilt University. His research on the design and evaluation of information systems artifacts has appeared in journals such as *ACM Computing Surveys*, *ACM Transactions on Database Systems*, *Communications of the ACM*, *Decision Sciences Journal*, *Decision Support Systems*, *IEEE Transactions on Knowledge and Data Engineering*, *Information Systems Research*, *Journal of MIS*, *Journal of Database Management*, and *MIS Quarterly*. He has served as the Editor-in-Chief of *ACM Computing Surveys*, as an associate editor for *MIS Quarterly*, and as a senior editor for *Information Systems Research*. He is currently senior editor emeritus for *Information Systems Research* and an associate editor for *Communications of the AIS*, *Decision Sciences Journal*, *Journal of Database Management*, *Information Systems and e-Business Management*, *The International Journal of Intelligent Information Technologies*, and *Information Systems Frontiers*.

# Appendix A

## Semantic Correctness Calculation ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬

Consider information request number three from the experimental task (Appendix E):

> Show the car rental agency, the date the car reservation was made, and the last name of the employee who made the car reservation for all car reservations made before 3/4/95.

A query that correctly satisfies this request for the state-based treatment (Figure 3) is

> select carrentalagency, carreserved, lastname
> from saleorder s
>   join employee e on e.employeeid=s.reservecaremployeeid
> where carreserved < '3/4/95'

A query that has several semantic errors appears below:

> select carreserved, lastname
> from saleorder s
>   join employee e on e.employeeid=s.instructorid
> where carreserved > '3/4/95'

This query has three errors. First, it selects car reservations made after 3/4/95 not before. Second, it omits one of the three required attributes (car rental agency). Third, the join uses an incorrect foreign key. Only the latter two are relevant to the experimental manipulations. Although the misuse of the greater-than sign is a semantic error, it is not involved in the mapping from subject's conceptualization of the domain to the related constructs in the data model and thus it is not considered to be relevant to the study.

Thus the semantic correctness of this query is evaluated as follows. There are seven required semantic elements, three attributes (carrentalagency, carreserved, lastname), two relations (saleorder, employee), and two keys (employeeid, reservecaremployeeid). The query identified five of the seven, so it has a semantic correctness score of 5/7 or 0.71.

# Appendix B

## Calculation of Mean Prediction Score ▬▬▬▬▬▬▬▬▬▬▬▬▬▬

Mean prediction score is a simple modification of the mean probability score (Yates 1990) to allow prediction of a continuous (rather than a dichotomous) variable. In calculating both, an arbitrary number of predictions are made, each with a stated confidence. For mean probability score, each prediction is either correct or incorrect, thus the outcome is coded as either one or zero. For mean prediction score, the outcome variable is the percent correct, and thus is coded anywhere from one to zero. The simple form of an individual prediction is as follows:

> Prediction Score = (prediction confidence – semantic correctness)$^2$

In the current study, the five-point Likert scale used to collect confidence was coded according to the following table:

| very high | 1.00 |
|-----------|------|
| high | 0.75 |
| medium | 0.50 |
| low | 0.25 |
| very low | 0.00 |

Thus, if a subject expressed high confidence in the accuracy of a query and identified all necessary semantic elements for a correct query, prediction score for that query would be calculated as follows:

$$\text{Prediction Score} = (0.75 - 1)^2 = 0.0625$$

If a subject expressed very high confidence in the accuracy of a query and identified 7 of 8 necessary semantic elements for a correct query, prediction score for that query would be calculated as follows:

$$\text{Prediction Score} = (1 - 0.875)^2 = 0.0156$$

Mean prediction score is simply the average of the prediction scores for each query.

# Appendix C

## Narrative Description of TechSupport Business Process ■■■■■■■■■■■■■■■■

TechSupport is a computer training company located in Salt Lake City, Utah. The company hires instructors and then contracts their services with other training companies. The database you will be working with is used to record the training activities of TechSupport.

When a client of TechSupport requests the services of an instructor, several steps need to be taken to ensure that the instructor will be at the client's location to deliver the requested training. Following is an example of the process used to organize the necessary arrangements for an instructor's trip to deliver a specific course.

Cathy answered the phone. It was Kevin from Executrain of Sacramento requesting an instructor to teach the 3-day Microsoft course on developing applications with Visual Basic. Executrain of Sacramento needed the class to be delivered in 3 weeks and Kevin hoped Stewart (the instructor who taught the class last time) would be available again. Cathy checked Stewart's schedule and indicated that he was available and that she would make the necessary arrangements to have Stewart in Sacramento to teach the class as requested. After the call, Cathy prepared the paperwork confirming the order and faxed it to Kevin, who signed it and faxed it back within an hour.

Now that the sale was complete, she began making the necessary arrangements. As usual, Gordon (her assistant) helped with the details. While Gordon called a travel agent to make arrangements for Stewart's airline ticket, Cathy placed an order with Microsoft to have the courseware delivered to Executrain of Sacramento's training facility. For the next 3 or 4 days, Cathy and Gordon worked to make sure that everything was in place for Stewart's trip (as they did for the many other trips which were currently planned for other instructors). Gordon called Stewart to make sure that he received his airline tickets and he also called Kevin at Executrain of Sacramento to make sure that Microsoft sent the correct manuals. Cathy made the arrangements for a hotel and rental car.

Everything went smoothly. Stewart taught the class in Sacramento and received high reviews from the students (as usual). As soon as the class was over, Cathy mailed an invoice for the price of instruction. When Stewart returned home, he submitted his expense report. Gordon forwarded the report to Executrain of Sacramento. After about three weeks a check came in the mail for the training and about a week later, a check to reimburse expenses. This is the process by which TechSupport conducts its business.

# Appendix D

## Relational Schemata Used in Experimental Treatments ■■■■■■■■■■

A combination of virtual tables and base tables implemented in Microsoft SQL Server comprise the logical implementation of each treatment. Virtual tables are defined using a create view statement. For example, the virtual table ConfirmSale, available only in the event-based treatment, is defined from the SaleOrder table as follows:

CREATE VIEW ConfirmSale as
SELECT ConfirmationSent, ConfirmationEmployeeID, SaleID
FROM SaleOrder;

Only subjects in the event-based treatment can reference this virtual table. Furthermore, as indicated below, subjects in the event-based treatment cannot reference the SaleOrder table as it is only available to subjects in the state-based treatment.

To illustrate how the same information is presented differently across treatments, consider car reservations (see Information Request 3 in Appendix E). The data items describing this phenomenon are the car rental agency, the date the car was reserved, and the employee who made the reservation (number and name). They are presented in each treatment as follows.

| State-Based Treatment | Event-Based Treatment | Artifact-Based Treatment |
|---|---|---|
| SaleOrder (table) | ReserveCar (virtual table) | Reservation (virtual table) |
|     CarRentalAgency |     CarRentalAgency |     CarRentalAgency |
|     CarReserved |     CarReserved |     CarReserved |
|     ReserveCarEmployeeID |     EmployeeID |     CarReservedEmployeeID |
| Employee (table) | Employee (table) | Employee (table) |
|     EmployeeID |     EmployeeID |     EmployeeID |
|     LastName |     LastName |     LastName |

The following relational schema details the underlying data structures used for each treatment. Primary keys are underlined. Table and virtual table names correspond directly to entity names in each treatment's conceptual-modeling script (E-R diagram). Each subject had access only to the tables and virtual tables that pertained to his or her treatment. Field names and descriptions were presented to subjects when they clicked on the corresponding entity in the E-R diagram as illustrated below for the Course table (available in all treatments).

| Course Table | 202 records |
|---|---|
| CertNoReq | Reference to a table not included in this database. |
| CourseID | Unique Identifier (Primary Key) |
| Days | The number of days it takes to teach the course |
| Name | The name of the Course |
| Price | Current list price for the course |
| ProductID | Reference to the Product table. Identifies the product the course is about |
| VendorCourseNumber | Identifier that the Vendor uses for this course |
| VendorID | Reference to the Vendor table. Identifies the company who produced the courseware |

## Tables and Virtual Tables Available in All Treatments

Course (CertNoReq, <u>CourseID</u>, Days, Name, Price, ProductID, VendorCourseNumber, VendorID) Table, 202 records
Customer (Address1, Address2, City, <u>CustomerID</u>, Fax, Name, Phone, PrincipleContact, RegionID, State, UsePerDiem, Zip) Table, 94 records
Employee (Address, City, <u>EmployeeID</u>, FirstName, LastName, Phone, State, Title, Zip) Table, 30 records
Hotel (Address1, Address2, City, <u>HotelID</u>, HotelName, Phone, State, Zip) Table, 85 records
PaymentType (PaymentDescription, <u>PaymentTypeID</u>) Table, 3 records
Product (GroupID, Name, <u>ProductID</u>, VendorID) Table, 54 records
Vendor (Name, Phone, <u>VendorID</u>) Table, 19 records

## Tables and Virtual Tables Available Only in the State-Based Treatment

Payment (Amount, Date, EmployeeID, PayTypeID, <u>PaymentID</u>, SaleID) Table, 797 records
SaleOrder (AirTicketsReceived, BeginDate, CarConfirmation, CarRentalAgency, CarReserved, CarType, ConfirmationEmployeeID, ConfirmationSent, Contact, CourseID, CourseVerifiedEmployeeID, CoursewareHandled, CoursewareVerified, CustomerID, ExpenseReportEmployeeID, ExpenseReportSent, HotelConfirmation, HotelID, HotelReserved, InstructorID, InvoiceEmployeeID, InvoiceSent, OrderTakingEmployeeID, OrderTicketsEmployeeID, Price, ReserveCarEmployeeID, ReserveHotelEmployeeID, <u>SaleDate</u>, SaleID, TicketsOrdered, TicketsRecdEmployeeID) Table, 610 records

## Tables and Virtual Tables Available Only in the Event-Based Treatment

ConfirmSale (ConfirmationSent, EmployeeID, <u>SaleID</u>) Virtual Table, 610 records
OrderTickets (EmployeeID, <u>SaleID</u>, TicketsOrdered) Virtual Table, 610 records
ReceiveOrder (BeginDate, Contact, CourseID, CustomerID, InstructorID, OrderTakingEmployeeID, Price, SaleDate, <u>SaleID</u>) Virtual Table, 610 records
ReceivePayment (Amount, Date, EmployeeID, PayTypeID, <u>ReceivePaymentID</u>, SaleID) Virtual Table, 797 records
ReceiveTickets (AirTicketsReceived, EmployeeID, <u>SaleID</u>) Virtual Table, 610 records
ReserveCar (CarConfirmation, CarRentalAgency, carReserved, CarType, EmployeeID, <u>SaleID</u>) Virtual Table, 610 records
ReserveHotel (EmployeeID, HotelConfirmation, HotelID, HotelReserved, <u>SaleID</u>) Virtual Table, 610 records
SendExpenseReport (EmployeeID, ExpenseReportSent, <u>SaleID</u>) Virtual Table, 610 records
SendInvoice (employeeID, InvoiceSent, <u>SaleID</u>) Virtual Table, 610 records
VerifyCourseware (CoursewareHandled, CoursewareVerified, EmployeeID, <u>SaleID</u>) Virtual Table, 610 records

## Tables and Virtual Tables Available Only in the Artifact-Based Treatment

Courseware (CoursewareHandled, CoursewareVerified, EmployeeID, <u>SaleID</u>) Virtual Table, 610 records
PaperWork (ConfirmationEmployeeID, ConfirmationSent, ExpenseReportEmployeeID, ExpenseReportSent, InvoiceEmployeeID, InvoiceSent, <u>SaleID</u>) Virtual Table, 610 records
Payment (Amount, Date, EmployeeID, PayTypeID, <u>PaymentID</u>, SaleID) Table, 797 records
Reservation (CarConfirmation, CarRentalAgency, CarReserved, CarReservedEmployeeID, CarType, HotelConfirmation, HotelID, HotelReserved, ReserveHotelEmployeeID, <u>SaleID</u>) Virtual Table, 610 records
SalesOrder (BeginDate, Contact, CourseID, CustomerID, InstructorID, OrderTakingEmployeeID, Price, SaleDate, <u>SaleID</u>) Virtual Table, 610 records
Ticket (AirTicketsReceived, OrderTicketsEmployeeID, <u>SaleID</u>, TicketsOrdered, TicketsRecdEmployeeID) Virtual Table, 610 records

# Appendix E

## Information Requests for Query Formulation Task ■■■■■■■■■■■■■■■■■■■■

1.  List the name and telephone number of all customers in the state of Colorado ("CO"). Sort the list alphabetically by name.
2.  List the Microsoft products for which TechSupport offers courses. Show product name and course name.
3.  Show the car rental Agency, the date the car reservation was made, and the last name of the employee who made the reservation for all car reservations made before 3/4/95.
4.  List the first and last names of employees who have made hotel reservations at the Holiday Inn. Also show the date the reservation was made and the city in which the particular Holiday Inn is located.
5.  Show the date on which the courseware for saleID 305 was verified.
6.  What is the total amount of money received from Executrain of Atlanta?
7.  What is the total price charged to Executrain of Atlanta for courses?

### Example Answers

| 1. | (All Treatments)   SELECT Name, Phone FROM Customer WHERE State = 'CO' ORDER BY Name |
|---|---|
| 2. | (All Treatments)   SELECT Product.Name, Course.Name FROM (Product JOIN Vendor ON Vendor.VendorID = Product.VendorID) JOIN course ON Course.ProductID = Product.ProductID WHERE Vendor.Name = 'Microsoft' |

| | Answers for the State-Based Treatment | Answers for the Event-Based Treatment | Answers for the Artifact-Based Treatment |
|---|---|---|---|
| 3. | SELECT   CarRentalAgency, CarReserved, LastName FROM SaleOrder s  JOIN Employee e ON e.EmployeeID = s.ReserveCarEmployeeID WHERE CarReserved < '3/4/95' | SELECT   CarRentalAgency, CarReserved, LastName FROM ReserveCar r  JOIN Employee e ON e.EmployeeID = r.EmployeeID WHERE CarReserved < '3/4/95' | SELECT   CarRentalAgency, CarReserved, LastName FROM Reservation r  JOIN Employee e ON e.EmployeeID = r.CarReservedEmployeeID WHERE CarReserved < '3/4/95' |
| 4. | SELECT   FirstName,LastName, HotelReserved, Hotel.City FROM SaleOrder s JOIN Employee e ON e.EmployeeID = s.ReserveHotelEmployeeID JOIN Hotel h ON h.HotelID = s.HotelID WHERE HotelName = 'Holiday Inn' | SELECT   FirstName,LastName, HotelReserved, Hotel.City FROM ReserveHotel r  JOIN Employee e ON e.EmployeeID =r.EmployeeID JOIN Hotel h ON h.HotelID = r.HotelID WHERE HotelName = 'Holiday Inn' | SELECT   FirstName,LastName, HotelReserved, Hotel.City FROM Reservation r  JOIN Employee e ON e.EmployeeID = r.ReserveHotelEmployeeID JOIN Hotel ON h.HotelID = r.HotelID WHERE HotelName = 'Holiday Inn' |
| 5. | SELECT   CoursewareVerified FROM SaleOrder WHERE SaleID = 305 | SELECT   CoursewareVerified FROM VerifyCourseware WHERE SaleID = 305 | SELECT   CoursewareVerified FROM Courseware WHERE SaleID = 305 |
| 6. | SELECT sum(Amount) FROM Payment p JOIN SaleOrder s ON p.SaleID = s.SaleID JOIN Customer c ON c.CustomerID = s.CustomerID WHERE c.Name = 'Executrain of Atlanta' | SELECT sum(Amount) FROM ReceivePayment p JOIN ReceiveOrder r ON p.SaleID = r.SaleID  JOIN customer c ON c.CustomerID = r.CustomerID WHERE c.Name = 'Executrain of Atlanta' | SELECT sum(Amount) FROM Payment p JOIN SalesOrder s ON p.SaleID = s.SaleID  JOIN Customer c ON c.CustomerID = s.CustomerID WHERE c.name = 'Executrain of Atlanta' |
| 7. | SELECT sum(price) FROM saleorder s JOIN customer c ON c.customerid = s.customerid  WHERE c.Name = 'Executrain of Atlanta' | SELECT sum(price) FROM receiveorder r JOIN customer c ON c.customerid = r.customerid WHERE c.Name = 'Executrain of Atlanta' | SELECT sum(price) FROM salesorder s JOIN customer c ON c.customerid = s.customerid WHERE c.Name = 'Executrain of Atlanta' |