

TSQL: A Design Approach

Richard Snodgrass

Department of Computer Science
University of Arizona
Tucson, AZ 85721
rts@cs.arizona.edu

May 8, 1992

I believe that many within the temporal database research community perceive that the time has come to consolidate approaches to temporal data models and calculus-based query languages, to achieve a consensus query language and associated data model upon which future research can be based. While some two dozen query language proposals exist, with a diversity of language and modelling constructs, common themes keep resurfacing. However, the community is quite fragmented, with each research project being based on a particular and different set of assumptions and approaches. Often these assumptions are not germane to the research *per se*, but are made simply because the research required a data model or query language with certain characteristics, with the particular one chosen rather arbitrarily. It would be better in such circumstances for research projects to choose the *same* language. Unfortunately, no existing language has attracted a following large enough to become the one of choice.

Gio Wiederhold of DARPA has long pressed for a consensus extension to SQL that could form a common core for future research. Let's term this extension the *Temporal Structured Query Language*, or TSQL (not to be confused with an existing language proposal of the same name). In this white paper I outline a proposal for a process by which a design for TSQL could be produced by the research community.

1 Scope

The scope of the TSQL language design should be restricted so that a coherent design is possible. In this section I propose aspects that should be included, and perhaps more importantly, those that should *not* be included.

- **TSQL is to be a *relational* query language.**

Given that SQL is “intergalactic dataspeak” (Mike Stonebraker’s term), TSQL should whenever possible be consistent with standard SQL, specifically, SQL89. It simply doesn’t make sense to base TSQL on competing (and arguably better) query languages such as Quel, Datalog, or Daplex. (Given that my language design work is based on Quel, I find it particularly painful that SQL has dominated over that superior language.) While I strongly agree that interesting research is possible and even desirable in extending the other languages to include temporal support, such extensions are necessarily outside of the scope of TSQL.

- **TSQL need not be consistent with existing standards.**

In general, TSQL need *not* be consistent with SQL2, which is in the final standardization process, nor SQL3, which is currently being designed. SQL2 contains severe flaws in its (minimal) handling of time-stamps, and SQL3 is a moving target which, in its present state, is regarded by many as a baroque design with a bewildering array of features. Consistency with the non-temporal aspects of existing standards for SQL, including SQL89, SQL2, and SQL3, is desirable if such consistency does not conflict with other goals.

- **TSQL will not be another standard.**

While the goal is a fully elaborated language design, there is no expectation that this design will be made into a standard. Of course, one hopes that our results would be acceptable to the standards bodies; at a minimum, our design should be communicated to these bodies. However, it is important to keep in focus the objective of the TSQL design: to provide a basis for future *research* in temporal databases. It also must be emphasized that TSQL should in no way limit or constrain future research in temporal databases, which should be free to adopt or propose whatever linguistic constructs are appropriate.

- **TSQL will not be an object-oriented query language.**

While temporal object-oriented query languages are being actively investigated, it would be distracting and counter-productive at this stage to attempt to merge the rather disparate approaches of object-oriented and relational languages while also addressing the temporal processing needs. Those involved in object-oriented language design are encouraged to produce, in parallel with this effort, a temporal object-oriented extension to SQL. At a later date, the two extensions could be merged.

- **TSQL should be comprehensive.**

TSQL should have constructs, extended in a natural fashion, that support all of the functionality of SQL, including update, aggregates, and schema specification and evolution. Consistent with the modifier “temporal”, TSQL should support both valid and transaction time.

- **The language design should include a formal semantics.**

Fortunately, there is a tradition of rigor in the temporal database community. The recent publication in TODS of a straightforward semantics of SQL will also help here.

- **The language will have an associated algebra.**

Such an algebra would demonstrate the existence of an executable equivalent to the declarative constructs in the language, and would suggest implementation strategies.

- **TSQL will be a *language design*.**

The TSQL design should not attempt to define storage structures, indexing structures, access methods, fourth-generation interfaces, support for distributed systems or heterogeneous databases, or optimization techniques. Such aspects, while important, are more properly the target of the research efforts that will utilize TSQL as a common substrate.

- **TSQL should reflect areas of convergence.**

The design of TSQL should avoid active areas of research where new results are generated frequently. Such areas include historical indeterminacy and temporal database design.

2 Language Design Process

It is in everyone's best interest to have as many participants in the design as possible. It would be wonderful to tap the extensive expertise available in the research community. On the other hand, the process must balance the desirability for input with the necessity of a design by a small number of designers, to avoid "design by committee" and all the difficulties such a design necessarily brings upon itself. Fortunately, there is a natural limiting mechanism available. Simply put, the design should be done by those researchers willing to expend the (considerable) effort to produce initial proposals and/or to modify designs in response to comments from a much larger community of evaluators.

Language designers will be self-selected persons who are willing to write *white papers* on some specific aspect of the design. White papers will include a survey of relevant research and a concrete proposal for some component of TSQL. Generally the proposal will include a formal syntax of the suggested constructs, an informal semantics (in prose) of these constructs, and, ideally, a formal semantics. These white papers should explicitly state the rationale behind important design decisions, to enable concrete discussion of the proposals.

Evaluators will be self-selected persons willing to comment in writing on a white paper. The comments will be collected, and addressed either by the author(s) of the initial white paper or by other designers willing to produce a new draft of the white paper.

The process will be iterative, and will converge when everyone is satisfied or exhausted. Meetings and workshops may be held to speed the design along.

Dissemination of the design can also be incremental, with consensus white papers being published in such outlets as *ACM SIGMOD Record* and *IEEE Data Engineering*. Conference panels might also be appropriate.

3 Tasks

Here I list a series of tasks that culminate in a fully elaborated language design. Each task has as its goal the production of and agreement on a white paper that addresses the indicated portion of the language.

3.1 Terminology

An agreed-upon set of concepts must be the first order of business. Fortunately, several researchers are working on exactly this issue in conjunction with a book being written on temporal databases.

It appears that convergence will be achieved soon.

3.2 Physical Time Line and Time-stamp Representation

Current DBMS's assume a time line starting at 1 A.D. or later and consisting of days or seconds, up to 9999 A.D.. One difficulty is that there are several definitions of *second* and of *day*. Another difficulty is that such a limited time line is of little use to many potential users of a temporal database, such as geologists, archæologists, anthropologists, and astronomers. Such a time line doesn't even include all of recorded history, and so doesn't fully support historians. Expanding the time line back to the creation of the universe (approximately 15 billion years ago), raises other definitional questions. For example, a solar year in the time of the dinosaurs was 400 days long. A year is difficult to define more than 6 billion years ago, before the earth was formed.

What is needed is an application-independent identification of one or more physical clocks that cover all of past time (15 billion years) and all of the foreseeable future. This definition of a physical time line should be convertible to other definitions that might be useful. A representation as a time-stamp data structure is also needed, with a precise semantics, i.e., a correspondence with a particular time of this physical clock for each valid bit pattern. Decisions need to be made about treating events as infinitely small points in time or as chronons of finite but nondecomposable length, closed or open representations for intervals, granularity, discrete versus continuous time, bounded versus unbounded time, and linear versus branching time.

3.3 User-defined Time Domain

In conventional as well as time-oriented databases, individual attributes can be associated with a temporal domain, termed *user-defined time*. Such a domain is supported by the DBMS in similar ways to other specialized domains, such as money, e.g., conversion to and from a string representation and the availability of comparison predicates. While SQL2 and DB2's SQL include two time-oriented attribute domains, *datetimes* and *intervals*, these language variants are limited to a single calendar, the Gregorian calendar, offer little or no support for anchored intervals, do not support languages other than English, and exhibit many problems with the semantics of arithmetic and boolean expressions. A proposal is needed that addresses these problems, while providing appropriate constructs for schema definition, time value input and output, predicates, arithmetic manipulation, and temporal functions.

3.4 Underlying Historical Relational Data Model

Determining the correct data model underlying TSQL will probably be the most difficult of all the tasks. Unfortunately, and not coincidentally, this task is a central one, on which most of the other tasks are predicated. To focus the design, I advocate that time be added to the data model in two separate steps, with the first to add valid time and the second to later add transaction time.

A proposal is needed that confronts the controversies currently raging in the research community, including 1NF versus \neg 1NF, temporally grouped versus temporally ungrouped, tuple time-stamped versus attribute value time-stamped, homogeneous versus non-homogeneous, events versus intervals, interpolated versus stepwise constant data, recurrent events, and whether keys should be required.

3.5 Benchmark Queries

A basis is needed on which to compare language proposals. This task involves informally defining an example schema containing several relations, populating this schema with example relation

instances, listing in English prose interesting queries on this schema, and displaying the results of these queries on the example instances.

3.6 Historical Selection and Projection

Historical selection is the analogue of conventional selection: the identification of tuples that satisfy some specified predicate, in this case a predicate on the time(s) the data elements (attribute values or tuples) were valid. One design issue is whether the **where** clause in SQL should be extended, or whether a new clause, such as the **when** clause in TQel, is preferred.

Historical projection is an analogue of conventional projection, where component(s) of tuples are retained, in this case, components of the time(s) the data elements were valid. One fundamental question is whether the derived intervals must be subsets of the underlying intervals. A second design issue is whether the target list in SQL should be extended, or whether a new clause, such as the **valid** clause in TQel, is preferred. The subject of temporal joins, such as time intersection, time union, and temporal outer joins, also needs to be addressed here.

3.7 Aggregates

Extension of the current SQL aggregates is required, along with the definition of new time-oriented aggregates (e.g., **first**), of temporal analogues of aggregate variants such as **unique** (e.g., moving window), and of order-dependent predicates that operate on groups of tuples.

3.8 Schema Specification and Evolution

SQL has a create table statement. This will need to be extended to allow specification of time-varying relations in addition to conventional relations. Other meta-data, such as the nature of interpolation to be imposed on continuous data represented discretely, must be included in the schema. Also, the schema of a relation may need to be changed to indicate a conversion from a time-varying relation, or vice versa. This will probably be a particularly easy extension to design.

3.9 Add Transaction Time to the Data Model

Since transaction time is orthogonal to valid time, the design process will be simplified if these two aspects are attacked separately. The hope is that once the impact of adding valid time to the language has been adequately considered, the incorporation of transaction time will be easier. Some feel that transaction time can be handled identically or almost identically to valid time; clearly if this is possible it will simplify this task considerably.

3.10 Schema Versioning

When schema evolution and support for transaction time are both present, a database may contain multiple versions of the schema, each in effect for disjoint intervals of transaction time. This aspect, while difficult to implement, probably has little impact on the language design.

3.11 Transaction Time Selection and Projection

At a minimum, these constructs should support rollback. A design decision is whether the historical selection and projection constructs should be extended, or whether different constructs, such as the **as of** clause in TQel, are needed.

3.12 Incorporate All SQL Constructs

To arrive at a comprehensive language definition, the interaction between proposed language constructs concerning time and all existing constructs, including modules, embedded SQL, views, and protection, needs to be examined in a systematic fashion.

3.13 Core Algebra

This task involves the design of representations for temporal relations (the objects in the algebra) and operators on these objects to support historical selection and projection. Issues including uni-sorted versus multi-sorted, algebraic equivalences, closure, snapshot reducibility, and update semantics should be considered.

3.14 Add Aggregates to the Algebra

Clearly the algebra should support all the aggregate variants present in the TSQL design.

3.15 Add Transaction Time to the Algebra

If schema versioning is to be supported in the algebra, this task must consider how algebra expressions are to be type checked in the presence of multiple schemas active at various transaction times.

4 Prerequisites

Dependencies between the tasks result in a partial order on their completion, as shown in Figure 1. These dependencies take the following considerations into account.

- The design of the time-stamp representation and of the user-defined time domain are independent of extensions of the underlying data model to incorporate valid or transaction time.
- The constructs for historical selection should be consistent with those used in expressions involving user-defined time.
- The constructs for historical selection and projection, unlike those for user-defined time, require a specified data model.
- While constructs for schema specification do not require the operations of historical selection and projection to be elaborated, the constructs for schema evolution will require new tuples to be generated that are consistent with a modified schema.
- Schema versioning is only possible if transaction time is supported.

While white papers that address only one task are best, those that address several tasks in concert will certainly still be welcome. Also, these dependencies are only suggestive; work can certainly proceed on multiple tasks concurrently.

5 Mechanics

The first step is for the research community to reach a consensus on the *process* of language design. Such a process has been outlined in this white paper. Either a modification of this paper, or a replacement with an alternate proposal, is a necessary pre-condition for subsequent effort.

I volunteer to serve as a clearinghouse for white papers on the tasks listed above. I will maintain a postal mailing list of people interested in receiving white papers. Additions to this list, as well as corrections and deletions, will be sincerely appreciated. Such modifications may be sent to me at either the postal or email addresses listed at the beginning of this paper. An email mailing list has also been set up, to enable interactive discussion of the TSQL design. To add or remove yourself from this mailing list, send mail to `tsql-request@cs.arizona.edu`. To send a message to all those on this mailing list, send mail to `tsql@cs.arizona.edu`. These last two addresses pertain to the email mailing list only.

Acknowledgements

The comments of many people who read a previous draft, especially those of Jim Clifford, Shashi Gadia, Sushil Jajodia, Christian Jensen, Sham Navathe, Arie Segev, and Abdullah Tansel, are appreciated. However, complaints and criticisms should be directed to the author.

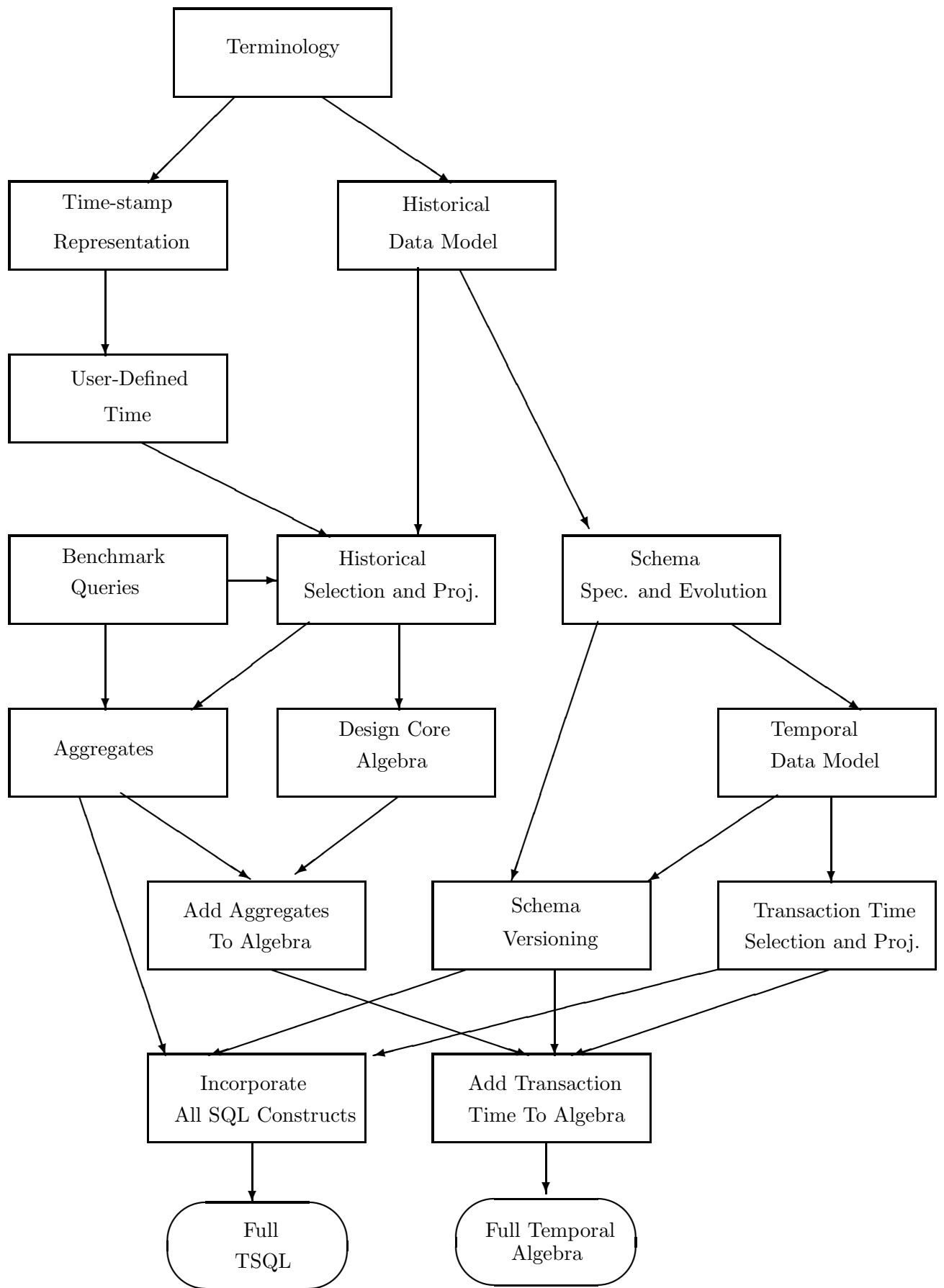


Figure 1: Task Dependencies