

TransactiveDB: Tapping into Collective Human Memories

Michele Catasta[†], Alberto Tonon^{*}, Djellel Eddine Difallah^{*}, Gianluca Demartini^{*},
Karl Aberer[†], and Philippe Cudre-Mauroux^{*}

[†]EPFL—Switzerland
(firstname.lastname)@epfl.ch

^{*}eXascale Infolab, University of Fribourg—Switzerland
(firstname.lastname)@unifr.ch

ABSTRACT

Database management systems (DBMSs) have been rapidly evolving in the recent years, exploring ways to store multi-structured data or to involve humans processes during query execution. In this paper, we outline a future avenue for DBMSs supporting transactive memory queries that can only be answered by a collection of individuals connected through a given interaction graph. We present TransactiveDB and its ecosystem, which allow users to pose transactive queries in order to reconstruct collective human memories. We describe a set of new transactive operators including `TUnion`, `TSelection`, `TJoin` and `TProjection`. We also describe how TransactiveDB leverages such transactive operators—by mixing query execution, social network analysis and human computation—in order to effectively and efficiently tap into the memories of all targeted users.

1. INTRODUCTION

Humans can store, process—and eventually forget—personal memories on their own, but also collectively; The notion of *transactive memory* [6] was established almost 30 years ago to denote the capacity of groups of individuals to collectively store and retrieve knowledge. Classical examples of transactive systems are older couples or families living together. Even if an individual in the group cannot remember a specific fact, he/she will often have a systematic way of retrieving the desired piece of information (e.g., by asking a specific individual from the group).

Recently, crowdsourcing has been used to complement classical database systems by leveraging human intelligence at scale [3]. While crowdsourcing DB systems can answer queries that purely automated systems cannot, they are still for the most part confined to relatively simple and generic tasks such as translating sentences from one language to another or tagging documents.

In this paper, we present our vision for TransactiveDB: a futuristic system able to leverage personal memories collectively and automatically in order to answer queries

whose results are not necessarily available in digital form, but are rather in the brains of relevant individuals. In Hippocampus [1], we presented a first example of transactive processing approach, which was able to reconstruct a shared memory more effectively than either machine-based or crowdsourcing-based approaches. We show how to generalize and systematize this approach in the following, and outline a generic architecture for a transactive DB system combining distributed query execution, human intelligence, and social networking to satisfy a novel class of information needs.

TransactiveDB works as a Peer Data Management System [4] in which each user represents an autonomous system combining relational tables (i.e., digital information) and human memories (i.e., personal information). Contrary to classical federated DB systems, the closed-world assumption is applied on the union of digital and personal information across the network of nodes. Both sources of information are fundamental to answer transactive queries expressing information needs that can only be satisfied by tapping into certain group memories, e.g., “Who is the person on the left of this picture that I took during the eXascale lab retreat in Anzere, Switzerland on Jan 30th, 2014?”

The rest of this paper is structured as follows. We start below by describing the rationale behind our system. Section 3 presents an overview of TransactiveDB architecture followed by the transactive operators that we introduce in section 4. We describe the social graph supporting TransactiveDB query execution in section 5. Section 6 outlines our query execution strategy. Finally, we lay out a research agenda for TransactiveDB and present our conclusions in sections 7 and 8, respectively.

2. SYSTEM RATIONALE

“What was the name of that amazing drink I ordered yesterday at the hipster bar?”

There are many reasons why a person might wish to ask such a question; more importantly for us is the fact that not everybody can provide an answer to it. That question actually translates into a transactive query that can neither be answered by querying the Web, nor by asking arbitrary Internet users via crowdsourcing. In order to give a precise answer to that query, one has to have interacted with the requester in the context of that query, meaning that one needs to have been with him at the hipster bar (direct interaction) or, at least, have heard about that event (indirect interaction). This question is an example of a *transactive point query*, a particular kind of transactive query supported by TransactiveDB that is processed iteratively and collaboratively by exploiting the interaction graph, that is, the graph

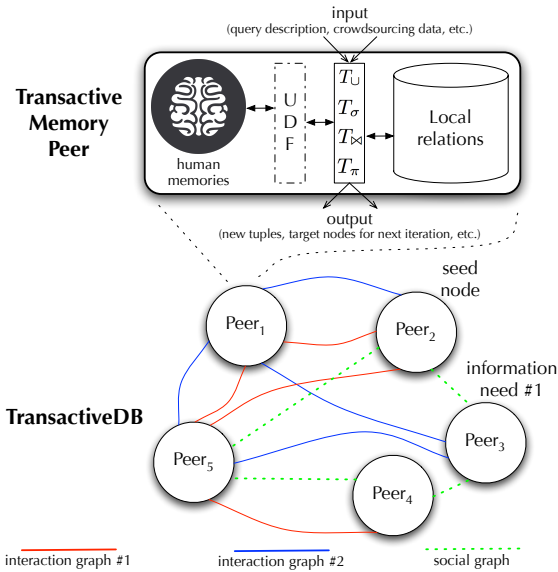


Figure 1: Architecture of TransactiveDB. In the depicted scenario, Peer₃ has the information need #1, which can be satisfied by the Transactive Memory peers belonging to the interaction graph #1. As such, Peer₃ exploits its social network connections to discover that Peer₂ can behave as the seed node for the query (i.e., Peer₂ can route the query to all the other peers involved in the relevant interaction graph).

representing all the relations among the people involved in the event the query refers to (a more detailed description of the interaction graph is given below, in Section 5). In that case, the point query targets one particular attribute (the name of the drink), but might also involve a second transactive element in order to determine the selection condition (i.e., to determine the name of the hipster bar).

In addition to point queries, TransactiveDB is also able to process further types of transactive queries such as transactive joins for matching entities, or transactive projections for completing missing data (see below Section 6). TransactiveDB could also be leveraged to reconstruct the transactive memory of an enterprise. In that case, employees could for example programmatically tap into the memories of their colleagues who participated in a given meeting in order to recover detailed events or decisions taken during a meeting.

A transactive DBMS sets itself apart from a crowd-augmented DBMS in the way participants are selected and incentivised, and in the way queries are collectively and iteratively executed. Where crowdsourcing employs a large number of anonymous, paid workers, our system leverages the acquaintances of the query requester to iteratively augment a shared corpus of information, i.e. the collective memory relating to the query. Subqueries are dispatched to the most suitable individuals depending on the interaction context (extending in that sense our previous push-crowdsourcing approach [2]). Nonetheless, we leverage classical crowdsourcing techniques to execute generic operators that do not require any specific knowledge, e.g., “Are these two photos depicting the same person?”.

3. OVERALL ARCHITECTURE

We envision a decentralized data management system that exposes memories of individuals or groups in order to answer transactive queries. A simplified architecture of TransactiveDB is depicted in Figure 1. The first key component of our system is the Transactive Memory Peer (a node hereafter), mainly composed of the memory of a person and a physical store onto which particular memories gets transcribed over time.

The physical store is a classical database system running for example in the cloud or on a personal device of the user. All the database system components are *transactive-enabled*, in the sense that they can execute transactive queries; in addition, these systems can also query the crowd via generic crowd operators (such as those defined in CrowdDB[3]). The memory transcription process can be triggered by two different types of events i) a voluntary act of documenting an experience, e.g., “record the following grocery list”, or ii) the reception of a query coming from the user or from a different node e.g.: “How much did last night meal at the restaurant cost?”.

The transcribed data is organized into relational tables that users create following a standard schema definition. With the proper security and privacy mechanisms (e.g., access control credentials), the data can be exposed and queried by trusted nodes in the system.

The second key element of our architecture is the *interaction graph* connecting the different nodes participating in query execution. The interaction graph is a subset of the underlying social network connecting the different end-users. The exact set of nodes and edges constituting the interaction graph is progressively elicited as the transactive query gets executed. Different queries can hence generate different interaction subgraphs (see Section 5).

4. TRANSACTIVE OPERATORS

TransactiveDB borrows from the standard relational algebra for basic operations. It uses operators defined by CrowdDB [3] for crowd-enabled queries. It also supports two new basic transactive operators, namely T_{Union} and $T_{Selection}$, and two derivative operators, namely T_{Join} and $T_{Projection}$, in order to handle transactive queries. It is up to the end-user providing the query to specify whether to use the crowdsourcing operations or the transactive memory ones in a declarative fashion. The set of operators we define below is not exhaustive, but regroups from our perspective all the key operators required to correctly express a transactive query.

T_{Union} . The T_{Union} operator, denoted by T_U , takes a relation R as input and returns a new relation R' containing all the tuples in R as well as new tuples retrieved transactively from other nodes. The experiment described in [1] is an application of the T_U operator and can be formalized by the following relational algebra operation:

$$all_participants \leftarrow T_U(initial_set),$$

where $initial_set = \{(name_1, sname_1, email_1), \dots, (name_n, sname_n, email_n)\}$ is the initial set of participants to the ISWC conference provided by the user who started the transactive memory experiment (notice that $initial_set$ may be an empty relation with a specified schema and, possibly, the starting seed).

TSelection. The TSelection operator, denoted by T_σ , takes as input a relation R and a predicate p , and exploits the transactive memory features of TransactiveDB to return a new relation R' composed of all the tuples of R that satisfy p . For example, to get all the participant of a conference that also attended the gala dinner, one can use the following operation, which leverages the transactive memory of the group to filter out participants that did not attend the dinner:

$$\text{hungry_participants} \leftarrow T_{\sigma_{\text{attended_dinner}=\text{true}}}(\text{all_participants}).$$

TJoin. The TJoin operator, denoted by T_{\bowtie} , is a ternary operator that takes as arguments two relations, R and S , and a predicate p . T_{\bowtie} leverages the transactive memory to execute a join between the tuples from two different relations satisfying the predicate. The transactive part of the operator selects which tuple of S should be associated with which tuples of R . This operator can be rewritten in different ways depending on the query and the instance data, opening the door to various query optimization strategies; when R and S are sufficiently small or when the interaction graph is sufficiently large, one can rewrite T_{\bowtie} as T_σ on $R \times S$ with p as predicate. In other cases, projecting on the join attributes and running TUnion queries to determine the values of those attributes separately before running the join might be more efficient.

TProjection. The TProjection operator, denoted by T_π , is a binary operator that takes as input a relation R and a set of attribute $A = \{a_1, \dots, a_n\}$. The output produced by $T_\pi(R, A)$ is the projection of the TUnion of R onto the attributes in A , that is,

$$T_\pi(R, A) = \pi_{a_1, \dots, a_n}(T_\cup(R)).$$

We note that the order of the selection and the TUnion is important: If the system first computes the selection, some contextual information required by the nodes to correctly process the transactive part of the query can be lost. For instance, with $T_\pi(\text{initial_set}, \{\text{email}\})$, where *initial_set* is as defined above, the system returns a list of emails from the participants of a conference, while with $T_\cup(\pi_{\text{email}}(\text{initial_set}))$ we obtain a generic set of e-mail addresses, since the peers might not have enough context to decide whose e-mails to contribute in that case.

5. LEVERAGING HUMAN MEMORIES THROUGH INTERACTIONS

TransactiveDB aims at reconstructing and gathering distributed information stored as human memories, be it digitally documented or dwelling inside the individuals' brains. Reconstructing such memories *a posteriori* requires some *a priori* exposure to the required information as well explicit social interactions to recompose the missing pieces. Social networks are today the canonical way of representing social interactions in a digital form. In fact, the existence of multiple social networks (e.g., professional, friendship, or celebrity networks) is a manifestation of the diverse social interactions human nurture. We define a *social graph* as the implicit or explicit graph encompassing the different interactions people have. Since a transactive query typically requires some very specific knowledge, it is only aimed at a subset of the *social graph*, called the interaction graph, see figure 2.

5.1 Interaction Graph

In our setting, collective memories are associated to a specific context where participants interact with each other, forming what we call an *Interaction Graph*. This notion is key to our system as it is created, expanded, and leveraged during query execution. To further illustrate this concept, we take the example of enumerating all participants of a conference (as we did in [1]); the interaction relates in that case to individuals participating to the particular conference; the interaction graph derived from this is a directed graph where the nodes correspond to attendees, and the edges are memories relating attendee A to attendee B . In a sense, *interaction graphs* are overlays sitting on top of *social graphs*.

5.2 Graph Creation

For queries relating to an interaction not yet observed in the system, TransactiveDB tries to discover the implicit interactions connecting the nodes by means of *connection elicitation query routing*. Practically, this requires recursively asking current nodes about further potential connections, and selectively routing the query to those nodes (old or new) that are the most susceptible of contributing new connections.

5.3 Graph Seed Selection

If the query requester did not take part in the interaction, the system engages in a seed discovery to bootstrap the interaction graph by iteratively exploring the social graph (e.g., similarly to the well-known Milgram experiment [5]). The seed is therefore defined as the first person whom the system identifies and who has taken part in the requested interaction. The quality of the seed plays an essential role in the efficiency of the subsequent transactional queries. For instance, selecting a potential “hub”, i.e., a person with many outgoing edges to further nodes in the interaction graph, is preferable to selecting more isolated nodes (in that sense, selecting the right seed participant for our conference attendance example might require ranking the potential seeds w.r.t. their connections in that field, their age, or their level of commitment for the conference.)

5.4 Leveraging the Interaction Graph

Interaction graphs are often created as new transactive queries arrive. As time goes by, however, queries relating to the same or similar interaction patterns may surface. TransactiveDB tries to reuse previously elicited interaction graphs for such queries. To follow up on our previous example, imagine a follow-up query asking for the “list of attendees of the benchmarking workshop that was co-located with the main conference”. In this case, the persons holding the information are already listed in the system thanks to the previous transactive query. This new query can hence leverage (at least in part) the interaction graph built for the preceding query directly.

6. QUERY EXECUTION

We give below an overview on how we envision queries to be executed through TransactiveDB. All operators are implemented as User Defined Functions (UDF) and exploit the interaction graph described in Section 5. Since the system may iteratively interact with its users in order to get missing information, a key factor we take into account is focus on humane-readable information. In that sense, tables, attributes and queries should all carry enough textual information to be self-explanatory.

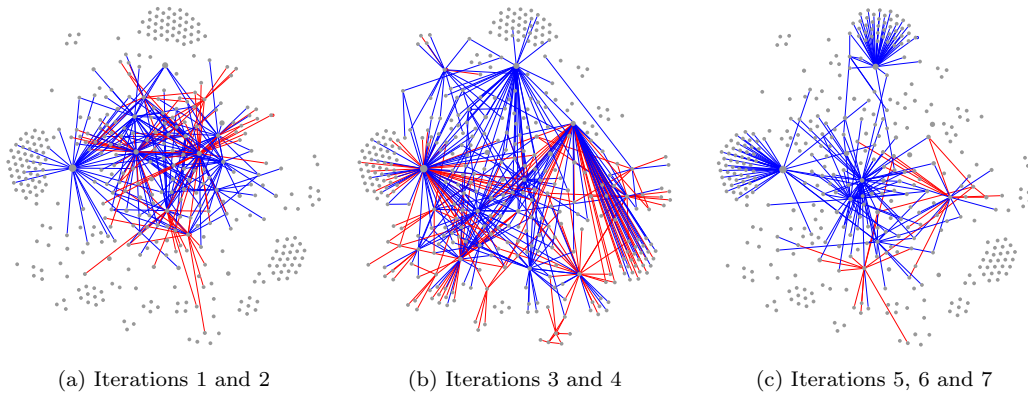


Figure 2: Interaction graphs while executing two transactive queries in parallel on the same set of peers. Each color encodes the new connections elicited by the different queries, while the visualization at different iterations shows how the contributions gradually move from central to peripheral peers.

The main steps of query execution are given in Algorithm 1. Query execution is iterative, and continues until a convergence criterion has been reached, or the budget used to operate the transactive memory system is exhausted.

At each step, the transactive memory peers are first selected and contacted based on the information contained in the interaction graph. Each Transactive Peer that the system reaches tries to retrieve the requested data using its local store. In case the information retrieved is deemed acceptable, the data is directly returned to the requester. Otherwise, the system generates a human-readable representation of the query, to which the local user can answer by filling out a Web-based form with information from his/her memory. At this step, missing data values can also be gathered through conventional crowdsourcing (e.g., emails of newly discovered participants). Finally, the results are post-processed and merged together and the peers used for next iteration are selected.

```

Data: empty or incomplete relation
Result: output relation after  $n$  iterations
initialization from seed node;
while converging or no more budget do
   $i$  = iteration number;
  contact Transactive peers reached at iteration  $i$ ;
  run Transactive Operators on each peer;
  collect results of iteration  $i$ ;
  if data missing in collected results then
    | fill gaps with crowdsourcing;
  end
  merge results with output relation;
  select target peers for iteration  $i+1$ ;
end

```

Algorithm 1: Steps of the Transactive Query execution

7. RESEARCH AGENDA

A number of research challenging relating to computer science and social sciences arose while designing TransactiveDB. One must tackle the following issues in our context: 1) social network analysis techniques should be developed in order to effectively and efficiently discover and determine the right nodes to query; 2) proper incentive mechanisms tailored to interaction networks should be devised to obtain information from the members of the transactive memory in a timely fashion. For instance, one could supplement our incentive scheme with ideas borrowed from collaboration platforms like Wikipedia, instead of relying on monetary incentives only; 3) malleable schema techniques coping with the heterogeneity of

the information handled by each individual should be designed; 4) agreement mechanisms and probabilistic metrics capturing the supposed validity of the records gathered through the transactive process should be used; 5) statistical tools should be developed in order to understand if the transactive process reached a point of convergence and if its execution should be stopped.

8. CONCLUSIONS

In this paper, we presented our vision for TransactiveDB a new DBMS capable of answering transactive memory queries. TransactiveDB fills a narrow, but in our opinion important, gap that neither search engines nor standard crowdsourcing technologies can fill and where the main source of information is the collective memory of a group of individuals connected through specific interactions. TransactiveDB features the usual relational algebra operators, the crowd operators of CrowdDB [3] and four new operators, namely TUnion, TSelection, TJoin, and TProjection, that allow the user to express transactive memory queries and to exploit the resulting graphs (“interaction graphs”) in order to discover which parts of the user’s social network are connected to the interactions processed by the system. We defined the key components of our envisioned system, and outlined a research agenda that we plan to further expand in the coming years in order to implement our vision.

9. REFERENCES

- [1] M. Catasta, A. Tonon, D. E. Difallah, G. Demartini, K. Aberer, and P. Cudré-Mauroux. Hippocampus: Answering Memory Queries using Transactive Search. 23rd International Conference on World Wide Web (WWW 2014), Web Science Track, 2014.
- [2] D. E. Difallah, G. Demartini, and P. Cudré-Mauroux. Pick-a-crowd: Tell Me What You Like, and I’ll Tell You What to Do. WWW ’13, pages 367–374, 2013.
- [3] M. J. Franklin, D. Kossmann, T. Kraska, S. Ramesh, and R. Xin. CrowdDB: Answering queries with crowdsourcing. SIGMOD ’11, pages 61–72, 2011.
- [4] A. Halevy, Z. Ives, J. Madhavan, P. Mork, D. Suciu, and I. Tatarinov. The piazza peer data management system. *Knowledge and Data Engineering, IEEE Transactions on*, 16(7):787–798, July 2004.
- [5] S. Milgram. The small world problem. *Psychology today*, 2(1):60–67, 1967.
- [6] D. M. Wegner. Transactive memory: A contemporary analysis of the group mind. In *Theories of group behavior*, pages 185–208. 1987.