

Hyperdatabases

Hans-Joachim Klemme, Böhmi, Torst Grabs,
Ulrich Rehm, Heiko Schuldt, Rog Weber
Swiss Federal Institute of Technology Zurich, Switzerland
{schek,rehm,grabs,boehmi,schuldt,weber}@inf.ethz.ch

Abstract

When relational database systems have been introduced twenty years ago, they were considered as infrastructure and a platform for application development. With today's information systems, the database system is no longer a management layer between the application and the data. Hyperdatabases become available to extend database concepts to high-level applications. A hyperdatabase manages distributed objects, software components, and workflow analogously to a database system that manages data and transactions. In short, hyperdatabases are called "higher order databases" which provide higher order data independence" e.g., immunity of applications against changes in the implementation of components and workload transparency. They will be the infrastructure for distributed information systems engineering of the future. In this abstract, we describe the current infrastructure and middle technology. This article elaborates on this vision. We will describe concrete projects: FT/HSZ, a PowerDB, a database cluster project, WebKow, an efficient document engine, and a build on top of database clusters. A further project studies transactional process management in a top of database transaction. Image similarity in multimedia components is another project. The hyperdatabase coordinates specialized components such as feature extraction and indexing services in a distributed environment.

1 Introduction Background

When relational database systems have been introduced twenty years ago, they have been considered as infrastructure and a platform for development of data-intensive applications. Standard database text lists the main motivation in the introduction. Independence was considered to be a breakthrough:

programmers were freed from low-level details, e.g., how to access shared data efficiently and correctly. It was even possible to have concurrent access. But now, the prerequisites for application development have changed dramatically. For instance, communication has become fairly cheap and the internet dominates modern information infrastructures. Consequently, the role of database concepts is re-evaluated and will be determined. Let us consider current topics, e.g., the management of software components and applications services, distributed client/middleware/server computing, application frameworks, enterprise resource planning systems, XML and e-commerce. Undoubtedly, the database system plays an important role. However, it has degenerated to a storage management layer, far away from the applications. More and more researchers are making these observations and start questioning the role of databases for future distributed WWW information systems engineering. They orient their research to make well-founded and established database concepts more applicable following the paradigm of [Regr]. Observations that current databases have been designed with the technology of twenty years ago. They state that the three-tier architecture, data at the bottom and application code always on the middle tier, is not the best. They state that databases are bloated by by relational features, by procedural stuff, and by warehouse features. They conclude that we should rethink everything. In a keynote speech, Brodie [Bro99] states that the database means nothing because of the inability to cope with the increase in data and transaction volumes. Heterogeneity hinders interoperability and accessibility. Architectural complexity is another issue which engineering is depositing with the solution for warehousing and mining.

The VLDB community [98] established a working group for application development in the database conference group to distinguish between two main directions: (1) database technology, and (2) infrastructure for information system development. According to the program committee, the VLDB (Future) more details and the interesting observations on evolution of database research are referred to documents such as [Si+96]. Even though the general awareness of the problem slowly proliferates, the database group [TZuri] has been involved in database technology, but also in the cooperation between databases and specialized application services for the "Cooperative Service Management of Open Systems" (COSMOS) research framework since 1990 [SSW90]. Taking these observations together, while keeping competence in database technology, starting at high levels where workflow and distributed component and object management takes place. The results will influence the next-generation platform for information-system development. The only hyperdatabase contribution of the TZuri database group is this vision.

In the following, we elaborate this vision and describe concrete projects. The presentation is similar to [ScW00] but it is incomplete with respect to project description [AHST97]. It is related to our emphasis on workflow and process management.

2 The Hyperdatabase Vision

2.1 Why Hyperdatabase?

First of all, the definition of Hyperdatabase (HDB) is a database of databases. An HDB administers objects that are composed of objects and transactions that are composed of transactions. Like hyper-matrices in numerical analysis, are matrices whose elements are matrices. An HDB database is a primitive of which are a database. Therefore, instead of talking about hyperdatabases, we may also talk about higher-order databases. We will use the terms synonymously in more general settings, say that an HDB administers distributed component network environment and provides kind of higher-order data independence". Remember that data independence was a major

future
proposed
(Epre
r
the
to
ch
ogy,
years.
roups,
y,
tion
rk
rich
escribe
e

breakthrough when relational databases were propagated. Now we must strive for immunity of application programs from change in storage and access structure, but also in this point there is an argument against changes in location, implementation, workload, the number of replicas, software components and their services.

What is the difference between a DB and an HDB? A nutshell way to say is: DB is a platform for clients concurrently accessing shared data. We need data definition, data manipulation, and transactional interface. The DB under the cover performs query optimization, correctness, parallel access, recovery, persistency, load balancing, admission control, availability.

Similarly, an HDB is a platform for clients concurrently accessing shared application services. As opposed to shared data in a DB, an HDB has shared components and services. An interface of an HDB needs components, service definition and description, service customization, transactional processes encompassing multiple service invocations. The HDB, under the cover, performs optimization, requests, routing, scheduling, parallelization, correctness of concurrent accesses, flexible failure treatment, providing guaranteed termination, availability, flexible recovery, and availability.

2.2 Hyperdatabase Middleware

The general architecture where HDB concepts are used is shown in Fig. 1. Client (white boxes) invokes application service (black boxes) server component that coordinates (large rectangle) application invocation. It may serve several coordinated component services. It may invoke several lower-level application services there. This continues along the invocation hierarchy until finally reach leaf nodes. The node performs data access, specialized computation (gray boxes). Note that a coordinator node plays a typical middleware dual role: they are servers for clients and clients to the server. The repeated record of complexity because servers are nodes. They have additional clients accessing distributed system concurrently to others coming from higher level coordinators. Seen from this perspective, the HDB concept helps manage tier architecture that exist and will further be realistic information system infrastructure.

Many products and standards propose architectures providing ready-to-use functionality. HDBs such as TP-monitors, evolving object monitors [BoK99], CORBA [CORBA], Enterprise

Members: Arun, Brodie, Carey, Gray, Ioannidis, Scheu, Wang, Widom.

Resource planning systems, federated databases, and warehouses. All provide the desirable feature of a distributed database. The DB provides a

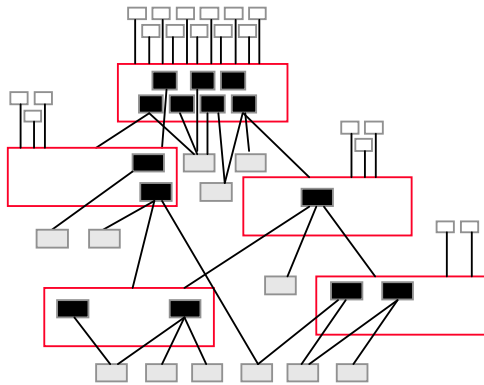


Figure 2.3: Repeated usage of a distributed architecture. Clients invoke application servers several coordinators.

all these approaches, incorporating the essential principles incorporated in database concepts but high level of abstraction with the records. This should allow for better understanding and convenient implementation. This objective of the HDB project is discussed in the following section. The current architecture sketch shows the opportunity for improvement.

2.3 TP Monitors

TP-Monitors are a narrow definition of a system for transaction processing. More generally, it is infrastructure for developing and running application services and components in a three-tier architecture. While TP-Monitors are definitely helpful in the use of many large-scale implementations, several room for extensions: a HDB would design for specification of distributed applications and automatic generation of program code as an extended workflow system. For example, HST97 further, a HDB would handle transaction layer above DB-transactions and transaction process management with failure treatment, availability guaranteed termination and "semantic" correctness. In current TP-monitor object monitors, components of a transactional TP-monitor that are forced to TP protocol. Coordination of all subtransaction performance level is a key only name. Finally, a HDB would handle optimal routing components and level "semantic" replication of the components.

data feature

2.4 Enterprise Resource Planning, SAP

An example consider the architecture of SAP system, a good representative ERP software. Clearly, observe the traditional DBMS for level Application node in the middle tier called application servers and presentation stays utmost level. The architecture of a HDB closely inspection of the middle layer the mapping of the DB to a AP clear that a HDB would be following: Most important, these should be a transaction on application objects above DB-transactions and transactional process management with semantic recovery and semantic correctness as well as the possibility of alternative execution as a separate transaction layer on a DB, a logical unit for which the sense HDB. However, application programmers must acquire locks and are responsible for correctness of parallel access and increase of failures. Also consistency of application server buffers is case of updates provided automatically. Second, a HDB would add flexible mapping many database storage manager in the current system on the database only, in order to allow a ability of one storage component is allowed, optimal decomposition and routing of client requests to managers that take place and level replication of whole storage component based on usage patterns achieved. Third, application servers should be transparent and would be under the responsibility of the application programmers, these with the current version of AP. A HDB would automatically optimize various cache and cache objects.

3 Hyperdatabase Projects

The following summarizes some of the projects at ETZ. Further, in more concrete explanation, what is already a HDB approach.

3. Composite Transactions and Transactional Process Management

Foundation. We have studied the problem of ensuring correctness of concurrent execution in a architecture. The figure 3.1 coordinates server component called coordinator in the following contains a full transaction manager for the client transactions

encompassing invocation of application over
 Every coordinator in the composite system performs
 transaction management on its own ensuring (local)
 correctness and (local) recovery. The problem is
 global correctness and global recovery. Insured,
 that each coordinator locally works correctly. In
 words we must know what additional ordering
 restrictions the coordinator must impose and
 over of "called" coordinator. What handshaking
 between coordinators is necessary to corre-
 ctly control concurrency in composite systems. In
 Fig. 1 we extensively studied this problem in
 past foundation of the project (e.g. We
 S92, WY93, Alo+97, Alo+99, Alo+99b),
 We S92, WY93, Alo+97, Alo+99, Alo+99b),
 studied practical protocols [SWS91, SSW95] and
 performed several evaluations [KaS96, NS96].

Transactional Process Management. In the
 transaction process management project we
 the conventional transaction model by grouping
 transaction into entities with high level semantic
 called *transactional processes*. These processes
 encompass flow control and the semantic
 elements. A side constraint for regular execution
 also possible specify alternative executions, wh
 applicable in failure. The phases have
 executed with processes in invocation of arbitr
 complex transaction provides services by
 of the low level. These transactions may dif
 terms of their termination behavior which indicate
 whether they do compensate and whether they
 to success after repeated invocation (retriability)
 [MRSK92, NBB94]. Given the termination behavior of
 single transaction invocations, the specificatio
 control flow and alternative single processes can
 proven correct. This is captured by the notion of
guaranteed termination, which generalizes the
 traditional notion of transaction atomicity. Having
 inherent properties, processes invoked termina
 well-defined state, correctly executing of
 several alternatives. A basic property of
 recovery (e.g. partial back and recovery combined
 alternative executions) another. The guaranteed
 termination will ensure several transac
 processes executed concurrently. The more com
 structure implies no complex dependencies compare
 with traditional transactions and be conside
 when transactional processes are scheduled [SAS99],
 thereby extending previous work on concurrency cont
 and recovery [Alo+97, WY93]. Similar approaches
 combining transaction management and process
 execution can be found in [WR92, CD96].
 Since we do not assume all applications managed by

coordinators depicted in Fig. 1. In a data
 systems, the transactional properties require for
 invocation may vary. In general, we have
 requirements we have analyzed in [SSA99], a
 transactional coordination agent (TCA) plugge
 application, thereby extending functionality adding
 transactional properties to service invocation (e.
 [NSSW94]). When considering, for instance,
 coordination processes in multimedia information
 system (which we discuss subsequently), these gen
 are extended to capture the workload of single
 components which are exploited for load-balan
 purposes.
 Work in transactions and transactional processes
 foundation and a basis for HDB transaction
 implementation. The principles have been applied
 for instance in the context of payment interaction in
 electronic commerce which are mapped processes,
 thereby making the execution guaranteed prov
 by transaction process management [SPS00].

3.2 PowerDB

The PowerDB project explores a distributed
 homogeneous component database (PC
 clusters) where every component is a complete DBMS
 with a Client access and a coordina
 with HDB (Fig. 2). We explore protocols for high-level
 transaction management in special consideration of
 replication. Replication of complete databases con
 tributes considerable speed-up in case of transacti
 ons.
 Due to the second layer transaction management we
 avoid a 2PC and avoid synchronous updates
 [GBS99, GBS00]. In addition query routing is
 detecting components that have the shortest respo
 nse time to queries that have been processed before
 [RBS00]. Replication can be restricted to
 certain parts of the database. We investigate metho
 ds that dynamically allow to choose components th

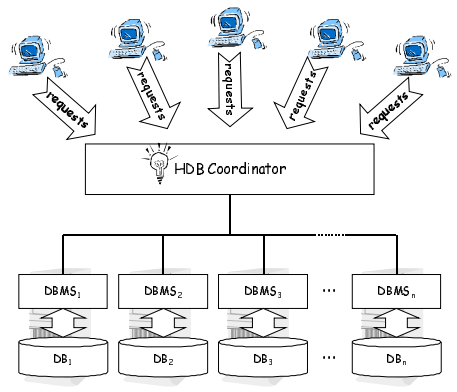


Fig. 2. PowerDB Architecture

The project "Routing OLAP Requests in a Cluster" and "Document Management on a DB Cluster" are described in more detail.

Routing OLAP Requests in a Cluster.

Analytical Processing queries (OLAP queries) refer to data warehouses, large collections of data, and accumulated data from operational databases. The complexity of users' OLAP queries is evaluated as follows.

The good retrieval performance is appropriate for physical data organization combined with query processing. If RBS00 has compared full replication and hybrid placement schemes combining partial replication with partitioning. Both approaches replicate data on some data nodes. Hence there are several components to be generated and evaluated. The Query Routing component decides on a DB node which component is to be evaluated. We have built a prototype system for PowerDB including among other features a secondary transaction management routing component providing different routing strategies and parallel distributed query processing based on [RB99]. Using this prototype, we have done performance evaluations with TPC-R benchmark comparing the different routing strategies: *simple round-robin-kind-of-strategy*, *random-first-come-first-server* or *Balance-the-Number-of-Queries* [CLL85], more sophisticated *affinity-based routing* strategy which assigns queries to a component (YCDI8 RBS00). The PowerDB architecture generally affects speedup of response time with increasing number of components (significant by hybrid placement scheme). Superduplication (showing even on thin speedup) has a significant effect on single fragments of partitioned relations. The smaller fragments of partitions which notably improve each hit rate per component. Besides we observed that OLAP queries evaluated concurrently typically obstruct each other leading to certain performance penalty. For example, the mean response time for executing queries concurrently on a node

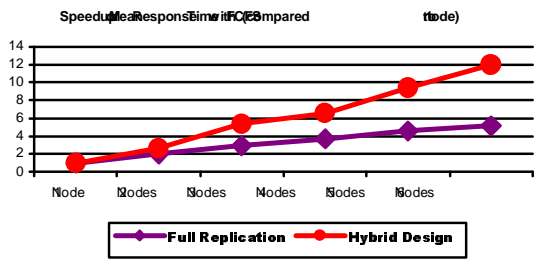


Figure 9: Query Routing Response Times

multiple programming views. The high availability of the cluster is maintained by using infinity-based query routing which improves the response time compared to CFS-routing by 30 percent on average. We are currently developing routing strategies which keep track of the component's actual state and route queries to the best-suited node.

Document Management on a DB Cluster.

[W3C98] proliferates a lot of information exchange for e-commerce. E-commerce applications require an efficient access to large collections of XML documents. The permanent availability of the documents that have been submitted to most today's document management systems is a substantial problem. This is due to the fact that the system has to be updated during off-line periods. XML databases where load is large are still prohibitive [FK99, ST+99]. The major drawback with this approach is that queries never operate on up-to-date data [KaS96, KaR96, BMV96]. Addressing this problem with systems that allow for concurrent retrieval of documents is a challenge. Both approaches are based on multi-level concurrency control [Wei99] derived from consistent scheduling with out sacrificing parallelism of the development. Our previous work is based on a database cluster and on forming the new document format XML have been addressed. Our current investigation on document management higher order objects and XML documents. The DB provides services for these higher order objects such as insertion, deletion, and update. Under the hood, the DB uses a service on a cluster of databases running off-the-shelf (similar to [ink96] [FG+97]). The DB additionally provides transaction management for document services. This leverages conventional database technology to document databases. The great advantage of this is the freedom to partition and replicate data and indexes to allocate them as many cluster components as necessary.

Decomposition and parallelization of requests together with composite transactions implemented at the coordination and the service layer of the DBs support short query response times and further allow for concurrent update and retrieval requests. Figure 10 shows the architecture of such a DB. We have implemented a prototypical document database for specific types of documents (i.e. for news group postings) [GBS99, GBS00]. A cluster of relational database systems stores the document text and the

inverted lists. Our experiments show that document insertions scale linearly and retrieval operations slightly

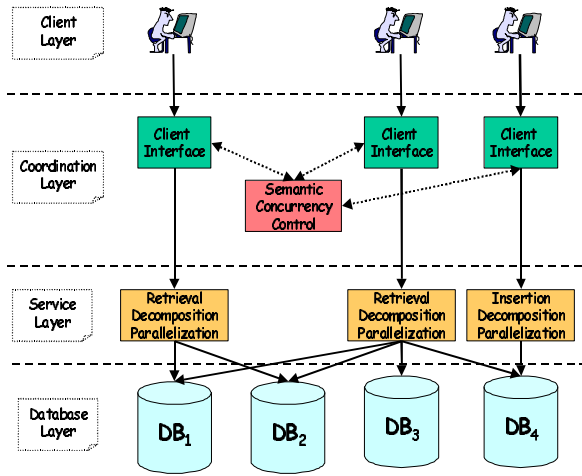


Figure 5: Document Management Hyperdatabase

less than for increasing cluster sizes. f. (Fig. 5). We observed that average response time in interactive for mixed workload of parallel insertions and retrieval streams reads directly from the database nodes. In our prototype system, our current work focuses on documents with arbitrary structure and restrictions on acceptable DBTs. An essential design issue is to reduce the amount of centralized processing in a distributed environment. This is important when the cluster size goes into some hundred database nodes that have no coordination. Our approach takes this into account as document specific functionality is distributed among the components. Only with information for semantic locking is necessarily kept in a centralized coordinator node. Results from GRS00 are encouraging. In our future work we want to address self-adapt hyperdatabase systems. The database developer system automatically partitions, replicates and materializes data for processing efficiently. At the service level, such a system replicates functionality among the coordinated nodes in order to deal with changing processing needs.

3.3 Image Search Management Systems.

Information systems for image collections consist of many specialized software components such as image servers, image processing, feature extraction, and indexing components. In such a setting given the possibly large number of components and the high workload imposed on them, location and workload

transparency is of great value and the HDB vision applies here. In more detail, the HDB Cluster, which contains many specialized software components installed as

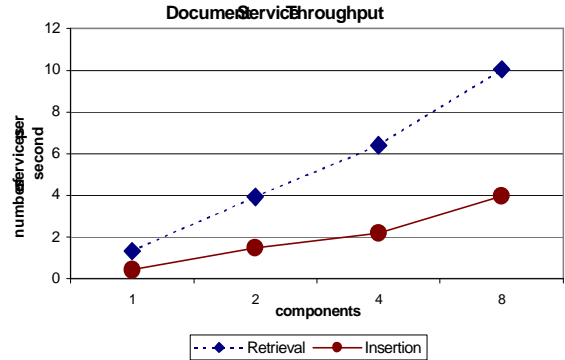


Figure 6: Document Service Scalability

necessary and coordinate them by a HDB. The coordination needs description of the software components including its actual load. Simple transaction processes for insertions, similarity and bulk load can be split into subtasks and optimally assigned to the components of the HDB. At any point in time, new components can be added to the cluster in order to improve response time. Interactive similarity retrieval is based on the VA-File, an efficient approximation of the inherently high-dimensional feature vector (WSB98, WeB00). In the following, we summarize the infrastructure aspects of parallelization necessary for interactive similarity and relevant feedback.

Coordination Management System Component

In order to coordinate the various components of the HDB, not only static information about the components is necessary, but also the dynamically changing state of every component at any point in time. For this purpose, we have introduced *coordination agents* [SSA99, Web+99] that are plugged into all components (fig. 6). A coordination agent observes the status and actions of its component and initiates processes that make sure the dependencies between components are properly maintained. Examples are processes for index maintenance, replication control, or consistency of metadata. The specification and execution of such processes is the main task of the HDB consistency agent. For instance, the image collection context repository agent of the repository initiates a **bulkload** process to extract the required features from the image and to set these features in a specialized index structure for similarity search.

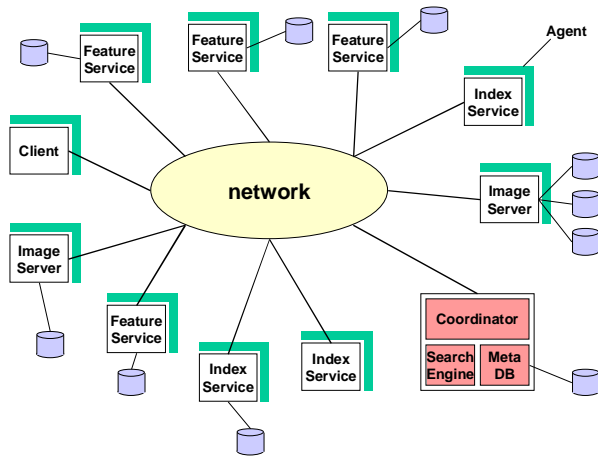


Fig. 6: Image System HDB

The main problem of bulkload processing is enormous cost of extracting features. For 100,000 images, pre-processing may take 50 days in a machine. Speedup the extraction HDB employs a large number of feature extraction components in parallel. Each component works on different images in the collection. The extraction of feature labels is long and dynamically update components in the system. After the registration of HDB, all the new components and out-requests them without the need of configuring or restarting the system. Even process can be started before adding new component can benefit from the additional component. This holds true for components retire from the system. In the coordination of HDB, both this, that components include of future processing and pending tasks are routed to available components. Load balancing is the key of HDB image database system. In bulkload process for example, pre-processing of individual depends on the size of the image. In the heterogeneous environment, to minimize the cost of process, must balance load over the available components equally. In the coordination of the extraction services periodically inform the HDB about the current load of the service. A course: time required working of all tasks assigned to the service. The HDB designs extraction of the service with the smallest current load and load low threshold. With this policy, can ensure that all components are almost equally loaded and the feature extraction components based on speed of current bulkload processes. More

details are found in [Boe+98, Web+99, WeS99, ScW00]. The same techniques for the process like similarity search are discussed in following.

Parallel Similarity Search and Relevant Feedback.

The rationale behind our activities in the field of similarity search is to provide search mechanisms that are user friendly and suitable for the user. Relevance and relevance feedback depends largely on the availability of a different feature order. Our relevance feedback mechanism adapts to the (subjective) relevance judgement of a user. Consequently, similarity search over image collection means nearest neighbour search (NN-search) in a high-dimensional feature space. A high-dimensional feature space is typically in the range of several hundred to thousands of dimensions. In this sense, similarity search over a large image collection requires considerable engineering effort in order to ensure an interactive response in this situation. HDB approach helps in following ways: the coordinator assigns feature extraction components and replicates indexing components as necessary, enabling high-level parallelization. NN-search in high-dimensional feature spaces is provably near the combinatorial complexity [WSB98]. In order to accelerate an avoidable linear scan, bitstring approximations of feature vectors. This structure is called the VA-File [WSB98, Web00]. NN-search in single VA-File follows a simple alpha-explicitly inspect a bitstring approximations. This gives a candidate for the nearest neighbour. In a second phase, check the candidate and determine the nearest neighbour. Note that this phase is well-suited for parallelization. We can partition the VA-File into arbitrary number of components. However, number of decisions rises as the number of issues whether to replicate or partition the approximation. The trade-off between flexibility and cost is more detailed flexibility means that the system can adequately react to changing behavior of component. For instance, the workload of a component may suddenly increase, or a component may become unavailable. The HDB is implemented as follows: the coordinator takes into account the current state of the components and finds an appropriate query evaluation strategy. This strategy remains transparent to the application. As the project described in this article, agents closely monitor the components and feed the coordinator with up-to-date information. Next, the workload the coordinator knows about the cache state of the components by approximating the

queries previously executed in a central account
 when assigning query components.
 What are the performance characteristics of this
 implementation of the HDB vision? A central
 observation is that a relatively small number of
 components is already sufficient for interactive
 answering time. We believe that a major
 since much research in similarity search did not
 any comparable solutions.
 Give such an efficient implementation of similar
 search, the natural next step is to allow the set
 formula to be informationally relevant and
 feedback mechanisms. I.e. the system collects user
 feedback and interprets it to refine the search
 to a certain step. This is a kind of a feedback
 to a statement in a similar query language and
 evaluate the query. A number of such mappings have
 been proposed in literature, e.g. [Roc71, RHOM98].
 A solution has been developed in a framework that is
 extensible and that allows for an integration of
 approaches that are around. A distinguishing component
 implements the relevance feedback functionality. It
 embeds the similarity search process actively
 easily thanks to our component-based approach. The
 specialized A-File component can easily be integra-
 ted and query answering time is again pleasingly low,
 even with complex relevant feedback queries.

4 Conclusions

In a first part we have represented the general
 hyperdatabase vision in a objective to provide
 higher-level platform for distributed application
 development. Hyperdatabase provides higher
 data independence analog to database systems
 and application programming will be lowered
 availability, correctness, well-defined termination,
 caching materialization, and replication of complete
 software components is merely an issue of the
 hyperdatabase system, taking off the responsibility
 from the programmer. Many successful tools
 realize this vision in the computer area
 some concrete projects of the THZ order exemplify
 some of these issues. We have summarized a transactional
 process management as a revolution from database
 transactional systems. We have shortly described PowerDB
 HDB Project with the main objective of demonstrating
 that coordinating and database systems in a distributed
 workstations pay off in terms of scalability and that
 there is no need for sacrificing correctness. The
 coordination of heterogeneous specialized components
 within a single system project. The A-
 File literature extraction is an example of very special
 components that nicely coordinate the HDB.

References

[Alo+97] Alonso G, Blot S, Fessler A, Schek, H.-J.:
*Correctness and Parallelism in Composite
 Systems.* In: Proceedings Principles
 of Database Systems (PODS'97), Tucson, Arizona,
 May 2-11 1997

[AHST97] Alonso G, Hager C, Sche H.-J., resc M.:
*Distributed Processing in Stand-alone Systems
 and Applications.* In: Proceedings International
 Conference on Very Large Databases (VLDB'97),
 August 1997, Athens, Greece

[Alo+99a] Alonso G, Fessler A, Pardon C, Sche H.-J.:
*Transaction Trac Fork in Composite
 Systems.* In: Proceedings International
 Theory of Database Systems (ICDT'99), Jerusalem, Israel, 10-12,
 (Beer, Bunema (Eds.), LNCS, vol. 1540,
 Springer-Verlag, 1999), pp. 150-168.

[Alo+99b] Alonso G, Fessler A, Pardon C, Sche H.-J.:
*Correctness in General Configurations of
 Transactional Components.* In: Proceedings
 ACM SIGMOD-SIGACT-SIGARTS Symp. on
 Principles of Database Systems (PODS'99),
 Philadelphia, Pennsylvania, May 1-Jun 4, ACM
 Press, New York, 1999), pp. 285-293.

[BGRS00] Böhm K, Grabst R, Röhm J, Schek H.-J.:
*Evaluating the Coordination Overhead of
 Synchronous Replication in a
 Databases.* To appear in Proceedings of the
 European Conference on Parallel Computing
 (Euro-Par'00), Munich, Germany, August 2000.

[BMV96] Barbara D, Mehrotra S, Valabhanan P.: *The
 Gold Text Indexing Engine.* In: Proceedings of the
 Twelfth International Conference on Data
 Engineering (ICDE'96), February 26 - March
 1, 1996, New Orleans, Louisiana, USA, pp. 172-179.

[Boe+98] Böhm K, Merj G, Schek H.-J., Weber,
 R.: *Metadata Management with HERMES
 Coordination Middleware.* ESPRIT project
 HERMES (no. 141), August 1998. Available at
<http://www-dbs.ethz.ch/~weber/paper/-HERMESmeta.ps>

[Boe00] Böhm K.: *Extending XEM Engine with
 Query Processing Capabilities.* In: IEEE
 Advances in Digital Libraries 2000, Bethesda,
 Maryland, USA.

[BoK99] Bouche K, Katz F.: *The Essential Guide to
 Object Monitors.* John Wiley & Sons, New York
 etd. 1999.

[Bro99] Brodie M.L.: *Query and Serendipitous
 Confluence of Economics, Business, and
 Collaborative Computing.* Proceedings of the
 International Conference on Data Engineering,
 Sydney, Australia, March 1999), pp. 3-3

[CD96] Chen Q, Dayal U.: *A Transactional Nested
 Process Management System.* In: Proceedings
 International Conference on Data Engineering
 (ICDE'96), New Orleans, Louisiana, February
 1996), pp. 66-73.

- [CLL85] Carey M.J., Livny M.L., H.: *Dynamic Task Allocation in Distributed Database Systems*. In: *Proceedings of the IEEE Conference on Distributed Computing Systems (ICDCS)*, Denver, Colorado, May 1985.
- [CORBA] <http://www.corba.org>
- [DFS99] Deutsch A., Fernandez M., Suciu D.: *Storing Semistructured Data with STORED*. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*, June 1999, Philadelphia, Pennsylvania, USA, pp. 431-442.
- [FG+97] Fogel G., Gribble S.D., Hawath B., Brewer E.A., Gauthier P.: *Cluster-Based Scalable Network Services*. In: *Proceedings of the ACM Symposium on Operating System Principles (SOSP'97)*, Malo, France, 1997, pp. 91.
- [FIK99] Florescu, D., Kossmann, D.: *Storing and Querying Multimedia in RDBMS*. In: *IEEE Database Engineering Bulletin*, 1999, 2(3), pp. 27-34.
- [GBS99] Grabt B., Böhm K., Schek H.-J.: *Document Engine in a Cluster*. In: *Proceedings of the Workshop on High Performance Transaction Systems (HPTS'99)*, Asilomar, California, Sept. 26-29, 1999.
- [GBS00] Grabt B., Böhm K., Schek H.-J.: *Parallel Document Engine in a Cluster of Databases -- Design, Implementation, and Experiences*. *Technical Report No. 40*, Dept. of Computer Science, ETH Zurich, April 2000.
- [GHOS96] Grahamellon B., Neill, E., Shashidhar, H.: *The Danger of Replication in Distributed Databases*. In: *Proceedings of the ACM SIGMOD Conference*, pp. 173-182, 1996.
- [ink96] Inktom Corp., *The Inktom Technology behind HotBot*, <http://www.inktomi.com/products-network/traffic/tech/clustered.html>, 1996.
- [KaR96] Kamath, M., Ramamritham, K.: *Efficient Transaction Support for Dynamic Information Retrieval Systems*. In: *Proceedings of the 9th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'96)*, Zurich, Switzerland, 1996.
- [KaS96] Kaufmann H., Schek H.-J.: *Extending TP-Monitoring to Intra-Transaction Parallelism*. In: *Proceedings of the International Conference on Distributed Database Systems (DDIS'96)*, Miami Beach, Florida, USA, Dec. 8-21, 1996, pp. 250-261.
- [MRSK92] Mehrotra S., Rastogi R., Silberschatz A., Korth H.: *A Transaction Model for Multidatabase Systems*. In: *Proceedings of the 2nd International Conference on Distributed Computing Systems (ICDCS'92)*, Yokohama, Japan, pp. 56-63.
- [NSSW94] Norrie M., Schaad W., Schek H.-J., Wunderli, M.: *CIM through Database Coordination*. In: *Proceedings of the International Conference on Data and Knowledge Systems for Manufacturing and Engineering (DKSME'94)*, Hong Kong, May 1994, pp. 668-673.
- [RB99] Röhms U., Böhm K.: *Working Together in Harmony: An Implementation of CORBA Object Query Services*. In: *Proceedings of the IEEE Conference on Data Engineering (ICDE)*, Sydney, Australia, March 1999.
- [RBS00] Röhms U., Böhm K., Schek H.-J.: *OLA Query Routing and Physical Design in a Database Cluster*. In: *Proceedings of the Conference on Extending Database Technology (EDBT'2000)*, Konstanz, Germany, March 27-31, 2000.
- [Regr] Seminar talk "A sketch of Regres" http://www.cs.berkeley.edu/~gribble/-summaries/talks_seminars/regres.html
- [RNS96] Rys, M., Norrie, M.C., Schek H.-J.: *Intra-Transaction Parallelism in Mapping an Object Model to a Relational Multi-Processor System*. In: *Proceedings of the 2nd VLDB Conf.*, Mumbai, Bombay, India, Sept. 4-9, 1996, pp. 460-471.
- [Roc71] Rocchio Jr., J.J.: *Relevance Feedback in Information Retrieval*. In: *MARIR Retrieval System Experiments in Automatic Document Processing*, Prentice Hall, 1970, Englewood Cliffs, New Jersey, USA, pp. 313-323.
- [RHOM98] Rui H., Huang T., Soragajam, Mehrotra S.: *Relevance Feedback for Interactive Content-Based Image Retrieval*. In: *IEEE Transactions on Circuits and Systems for Video Technology*, Special Issue on Segmentation, Description and Retrieval of Video Content, (5), 1998, pp. 644-655.
- [ScW00] Schek H.-J., Weber R.: *High Order Databases and Multimedia Information*. In: *Proceedings of the Swiss/Japan Seminar on Advanced Database and Multimedia*, Kyoto, Japan, Feb. 2000.
- [SAS99] Schuldt, Alonso G., Schek H.-J.: *Concurrency Control in Recovery of Transactional Process Management*. In: *Proceedings of the 18th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS'99)*, Philadelphia, Pennsylvania, May 31 - June 2, 1999, New York, pp. 316-326.
- [SPS00] Schuldt, Popovic, Schek H.-J.: *Automatic Generation of Reliable-Commerce Payment Processes*. In: *Proceedings of the International Conference on Web Information Systems Engineering (WISE'00)*, Hong Kong, China, June 2000.
- [SSA99] Schuldt, H., Schek, H.-J., Alonso, G.: *Transactional Coordination Agents for Composite Systems*. In: *Proceedings of the International Database Engineering and Applications Symposium (IDEAS'99)*, Montréal, Canada, August 1999, pp. 321-331.
- [Si+96] Silberschatz A., et al.: *Strategic Directions in Database Systems*. *Breaking the Box*. ACM Computing Surveys, Vol. 28, No. 4, Dec. 1996, pp. 764-778.

- [SSW90] Schek H.-J., Schödl M., Weikum G.: *From the KERNEL to the COSMOS: The Database Research Group at ETH Zurich*. Tech Report No. 3, Dept. of Computer Science, ETH Zurich, July 1990.
- [SSW95] Schaad, W., Schek, H.-J., Weikum, G.: *Implementation and Performance of Multi-level Transaction Management in a Multidatabase Environment*. In: *15th Workshop on Research Issues on Data Engineering Distributed Object Management (RIDE-DOM'95)*, Taipei, Taiwan, March 1995, pp. 8-115.
- [ST+99] Shanmugasundaram, T., Fu, K. H., Zhang, C., DeWitt, D., Naughton, J.: *Relational Databases for Querying XML Documents: Limitations and Opportunities*. In: *Proceedings of the 25th International Conference on Very Large Databases (VLDB '99)*, Sept. 7-10, 1999, Edinburgh, Scotland, pp. 302-314.
- [SWS91] Schek H.-J., Weikum G., Schaad W.: *Multi-Level Transaction Approaches for Federated DBMS Transaction Management*. In: *Proceedings of the First International Workshop on Interoperability in Multidatabase Systems (IMS'91)*, Kyoto, April 1991.
- [SWY93] Schek H.-J., Weikum G., Yeh S.-L.: *Towards a Unified Theory of Concurrency Control and Recovery*. In: *12th ACM SIGACT-SIGMOD-SIGAR Symposium on Principles of Database Systems (PODS)*, Washington DC, May 1993 (Appeared in *Tech. Report 93-19*, Dept. of Computer Science, ETH Zurich, Dec 1992).
- [VD98] Vogel, W., Dumitriu, D.: *The Design and Architecture of Microsoft Cluster Service: A Practical Approach to High-Availability and Scalability*. FTCS 98.
- [VLDB] <http://www.vldb.org/>
- [W3C98] The World Wide Web Consortium: *Extensible Markup Language (XML) 1.0 - W3C Recommendation 10-February-1998*. Available at: <http://www.w3.org/TR/1998/REC-xml-19980210>
- [Web+99] Weber R., Bolliger I., Grosz S., Schek, H.-J.: *Architecture of Network Image Search and Retrieval System*. In: *Proceedings of the Conference on Information Knowledge Management (CIKM'99)*, Kansas City, Missouri, No. 2, ACM Press, New York, 1999, pp. 430-441.
- [WeB00] Weber R., Böhm K.: *Trading Quality of Time with Nearest-Neighbor Search*. In: *Proceedings of the Conference on Extending Database Technology (EDBT 2000)*, Konstanz, Germany, pp. 1-35.
- [Wei91] Weikum G.: *Principles and Realization Strategies of Multi-Level Transaction Management*. In: *ACM Transactions on Database Systems*, vol. 6, No. 1, March 1991, pp. 32-180.
- [WeS92] Weikum, G., Schek, H.-J.: *Concepts and Application of Multilevel Transactions in Open Nested Transactions*. In: *Elmagarmy (Ed.), Database Transaction Models for Advanced Application*, Morgan Kaufmann, San Mateo, CA, 1992.
- [WeS99] Weber R., Schek H.-J.: *Distributed Image-Database Architecture for Efficient Insertion and Retrieval*. In: *Proceedings of the Workshop on Multimedia Information Systems (MIS'99)*, Indian Wells, California, Oct. 21-23, Golubchikov M. Tsotras, Eds., pp. 8-55.
- [WR92] Wächter, R., Reuter A.: *The Contraction Model*, chapter 19-26. In: *Elmagarmy (Ed.), Database Transaction Models for Advanced Application*, Morgan Kaufmann, San Mateo, CA, 1992.
- [WSB98] Weber R., Schek H.-J., Bloß G.: *Quantitative Analysis and Performance Study of Similarity-Search Methods in High-Dimensional Spaces*. In: *Proceedings of the Conference on Very Large Databases (VLDB '98)*, New York, USA, Aug. 24-27, 1998.
- [YCDI87] Yeh, S.-L.: *Analysis of Final Routing in Multi-System Data Sharing*. Performance Evaluation, 7:87-109, 1987.
- [ZNB94] Zhang, N., Nodim, M., Bhargava, A., Bukhr, S.: *Ensuring Relaxed Atomicity for Flexible Transaction Multidatabase Systems*. In: *Proceedings of the SIGMOD International Conference on Management of Databases (SIGMOD'94)*, Minneapolis, Minnesota, May 1994, pp. 67-78.