

PicoDBMS: Scaling down database techniques for the smartcard

Philippe Pucheral¹, Luc Bouganim¹, Patrick Valduriez², Christophe Bobineau¹

¹ University of Versailles, PRiSM Laboratory, Versailles, France;

E-mail: {philippe.pucheral;luc.bouganim;christophe.bobineau}@prism.uvsq.fr

² University Paris 6, LIP6 Laboratory, Paris, France; E-mail: patrick.valduriez@lip6.fr

Edited by A. El Abbadi, G. Schlageter, K.-Y. Whang. Received: 15 October 2000 / Accepted: 15 April 2001

Published online: 23 July 2001 – © Springer-Verlag 2001

Abstract. Smartcards are the most secure portable computing device today. They have been used successfully in applications involving money, and proprietary and personal data (such as banking, healthcare, insurance, etc.). As smartcards get more powerful (with 32-bit CPU and more than 1 MB of stable memory in the next versions) and become multi-application, the need for database management arises. However, smartcards have severe hardware limitations (very slow write, very little RAM, constrained stable memory, no autonomy, etc.) which make traditional database technology irrelevant. The major problem is scaling down database techniques so they perform well under these limitations. In this paper, we give an in-depth analysis of this problem and propose a PicoDBMS solution based on highly compact data structures, query execution without RAM, and specific techniques for atomicity and durability. We show the effectiveness of our techniques through performance evaluation.

Key words: Smartcard applications – PicoDBMS – Storage model – Execution model – Query optimization – Atomicity – Durability

1 Introduction

Smartcards are the most secure portable computing device today. The first smartcard was developed by Bull for the French banking system in the 1980s to significantly reduce the losses associated with magnetic stripe credit card fraud. Since then, smartcards have been used successfully around the world in various applications involving money, proprietary data, and personal data (such as banking, pay-TV or GSM subscriber identification, loyalty, healthcare, insurance, etc.). While today's smartcards handle a single issuer-dependent application, the trend is toward multi-application smartcards¹. Standards for multi-application support, like the JavaCard [36] and Microsoft's SmartCard for Windows [26], ensure that the card be universally accepted and be able to interact with several

service providers. This should make smartcards one of the world's highest-volume markets for semiconductors [14].

As smartcards become more and more versatile, multi-application, and powerful (32-bit processor, more than 1 MB of stable storage), the need for database techniques arises. Let us consider a health card storing a complete medical folder including the holder's doctors, blood type, allergies, prescriptions, etc. The volume of data can be important and the queries fairly complex (select, join, aggregate). Sophisticated access rights management using views and aggregate functions are required to preserve the holder's data privacy. Transaction atomicity and durability are also needed to enforce data consistency. More generally, database management helps to separate data management code from application code, thereby simplifying and making application code smaller. Finally, new applications can be envisioned, like computing statistics on a large number of cards, in an asynchronous and distributed way. Supporting database management on the card itself rather than on an external device is the only way to achieve very high security, high availability (anywhere, anytime, on any terminal), and acceptable performance.

However, smartcards have severe hardware limitations which stem from the obvious constraints of small size (to fit on a flexible plastic card and to increase hardware security) and low cost (to be sold in large volumes). Today's microcontrollers contain a CPU, memory – including about 96 kB of ROM, 4 kB of RAM, and up to 128 kB of stable storage like EEPROM – and security modules [39]. EEPROM is used to store persistent information; it has very fast read time (60–100 ns) comparable to old-fashion RAM but very slow write time (more than 1 ms/word). Following Moore's law for processor and memory capacities, smartcards will get rapidly more powerful. Existing prototypes, like Gemplus's Pinocchio card [16], bypass the current memory bottleneck by connecting an additional chip of 2 MB of Flash memory to the microcontroller. Although a significant improvement over today's cards, this is still very restricted compared to other portable, less secure, devices such as Personal Digital Assistants (PDA). Furthermore, smartcards are not autonomous, i.e., have no independent power supply, thereby precluding asynchronous and disconnected processing.

¹ Everyone would probably enjoy carrying far fewer cards.

These limitations (tiny RAM, little stable storage, very costly write, and lack of autonomy) make traditional database techniques irrelevant. Typically, traditional DBMS exploit significant amounts of RAM and use caching and asynchronous I/Os to reduce disk access overhead as much as possible. With the extreme constraints of the smartcard, the major problem is scaling down database techniques. While there has been much excellent work on scaling up to deal with very large databases, e.g., using parallelism, scaling down has not received much attention by the database research community. However, scaling down in general is becoming very important for commodity computing and is quite difficult [18].

Some DBMS designs have addressed the problem of scaling down. Light versions of popular DBMS like Sybase Adaptive Server Anywhere [37], Oracle 8i Lite [30] or DB2 Everywhere [20] have been primarily designed for portable computers and PDA. They have a small footprint which they obtain by simplifying and componentizing the DBMS code. However, they use relatively high RAM and stable memory and do not address the more severe limitations of smartcards. ISOL's SQLJava Machine DBMS [13] is the first attempt towards a smartcard DBMS while SCQL [24], the standard for smartcard database language, emerges. While both designs are limited to single select, they exemplify the strong interest for dedicated smartcard DBMS.

In this paper, we address the problem of scaling down database techniques and propose the design of what we call a PicoDBMS. This work is done in the context of a new project with Bull Smart Cards and Terminals. The design has been made with smartcard applications in mind, but its scope extends as well to any ultra-light computer device based on a secured monolithic chip. This paper makes the following contributions:

- We analyze the requirements for a PicoDBMS based on a typical healthcare application and justify its minimal functionality.
- We give an in-depth analysis of the problem by considering the smartcard hardware trends and derive design principles for a PicoDBMS.
- We propose a new pointer-based storage model that integrates data and indices in a unique compact data structure.
- We propose query execution techniques which handle complex query plans (including joins and aggregates) with no RAM consumption.
- We propose transaction techniques for atomicity and durability that reduce the logging cost to its lowest bound and enable a smartcard to participate in distributed transactions.
- We show the effectiveness of each technique through performance evaluation.

This paper is an extended version of [7]. In particular, the section on transaction management is new. The paper is organized as follows. Section 2 illustrates the use of take-away databases in various classes of smartcard applications and presents in more detail the requirements of the health card application. Section 3 analyzes the smartcard hardware constraints and gives the problem definition. Sections 4–6 present and assess the PicoDBMS' storage model, query execution model, and transaction model, respectively. Section 7 concludes.

2 Smartcard applications

In this section, we discuss the major classes of emerging smartcard applications and their database requirements. Then, we illustrate these requirements in further detail with the health card application, which we will use as reference example in the rest of the paper.

2.1 Database management requirements

Table 1 summarizes the database management requirements of the following typical classes of smartcard applications:

- *Money and identification*: examples of such applications are credit cards, e-purse, SIM for GSM, phone cards, transportation cards. They are representative of today's applications, with very few data (typically the holder's identifier and some status information). Querying is not a concern and access rights are irrelevant since cards are protected by PIN-codes. Their unique database management requirement is update atomicity.
- *Downloadable databases*: these are predefined packages of confidential data (e.g., diplomatic, military or business information) that can be downloaded on the card – for example, before traveling – and be accessed from any terminal. Data availability and security are the major concerns here. The volume of data can be important and the queries complex. The data are typically read-only.
- *User environment*: the objective is to store in a smartcard an extended profile of the card's holder including, among others, data regarding the computing environment (PC's configuration, passwords, cookies, bookmarks, software licenses, etc.), an address book as well as an agenda. The user environment can thus be dynamically recovered from the profile on any terminal. Queries remain simple, as data are not related. However, some of the data are highly private and must be protected by sophisticated access rights (e.g., the card's holder may want to share a subset of her/his address book or bookmark list with a subset of persons). Transaction atomicity and durability are also required.
- *Personal folders*: personal folders may be of a different nature: scholastic, healthcare, car maintenance history, loyalty. They roughly share the same requirements, which we illustrate next with the healthcare example. Note that queries involving data issued from different folders can make sense. For instance, one may be interested in discovering associations between some disease and the scholastic level of the card holder. This raises the interesting issue of maintaining statistics on a population of cards or mining their content asynchronously.

2.2 The health card application

The health card is very representative of personal folder applications and has strong database requirements. Several countries (France, Germany, USA, Russia, Korea, etc.) are developing healthcare applications on smartcards [11]. The initial idea was to give to each citizen a smartcard containing her/his identification and insurance data. As smartcard storage capacity increases, the information stored in the card can be

Table 1. Typical applications' profiles

Applications	Volume	Select/project	Join	Group by / Distinct	Access rights / views	Atomicity	Durability	Statistics
Money & identification	tiny					✓		
Downloadable DB	high	✓	✓	✓				
User environment	medium	✓			✓	✓	✓	
Personal folder	high	✓	✓	✓	✓	✓	✓	✓

extended to the holder's doctors, emergency data (blood type, allergies, vaccination, etc.), surgical operations, prescriptions, insurance data and even links to heavier data (e.g., X-ray examination, scanner images, etc.) stored on hospital servers. Different users may query, modify, and create data in the holder's folder: the doctors who consult the patient's past records and prescribe drugs, the surgeons who perform exams and operations, the pharmacists who deliver drugs, the insurance agents who refund the patient, public organizations which maintain statistics or study the impact of drugs correlation in population samples, and finally the holder her/himself.

We can easily observe that: (i) the amount of data is significant (more in terms of cardinality than in terms of volume because most data can be encoded); (ii) queries can be rather complex (e.g., a doctor asks for the last antibiotics prescribed to the patient); (iii) sophisticated access rights management using views and aggregate functions are highly required (e.g., a statistical organization may access aggregate values only but not the raw data); (iv) atomicity must be preserved (e.g., when the pharmacist delivers drugs); and (v) durability is mandatory, without compromising data privacy (logged data stored outside the card must be protected).

One may wonder whether the holder's health data ought to be stored in a smartcard or in a centralized database. The benefit of distributing the healthcare database on smartcards is threefold. First, health data must be made highly available (anywhere, anytime, on any terminal, and without requiring a network connection). Second, storing sensitive data on a centralized server may damage privacy. Third, maintaining a centralized database is fairly complex due to the variety of data sources. Assuming the health data is stored in the smartcard, the next question is why the aforementioned database capabilities need to be hosted in the smartcard rather than the terminals. The answer is again availability (the data must be exploited on any terminal) and privacy. Regarding privacy, since the data must be confined in the chip, so must the query engine and the view manager. As the smartcard is the unique trusted part of the system, access rights and transaction management cannot be delegated to an untrusted terminal.

3 Problem formulation

In this section, we make clear the smartcard constraints in order to derive design rules for the PicoDBMS and state the problem. Our analysis is based on the characteristics of both

existing smartcard products and current prototypes [16, 39], and thus, should be valid for a while. We also discuss how the main constraints of the smartcard will evolve in a near future.

3.1 Smartcard constraints

Current smartcards include in a monolithic chip, a 32 bits RISC processor at about 30 MIPS, memory modules (of about 96 kB of ROM, 4 kB of static RAM, and 128 kB of EEPROM), security components (to prevent tampering), and take their electrical energy from the terminal [39]. ROM is used to store the operating system, the JavaCard virtual machine, fixed data, and standard routines. RAM is used as working memory for maintaining an execution stack and calculating results. EEPROM is used to store persistent information. EEPROM has very fast read time (60–100 ns/word) comparable to old-fashion RAM, but a dramatically slow write time (more than 1 ms/word).

The main constraints of current smartcards are therefore: (i) the very limited storage capacity; (ii) the very slow write time in EEPROM; (iii) the extremely reduced size of the RAM; (iv) the lack of autonomy; and (v) a high security level that must be preserved in all situations. These constraints strongly distinguish smartcards from any other computing devices, including lightweight computers like PDA.

Let us now consider how hardware advances can impact on these constraints, in particular, memory size. Current smartcards rely on a well-established and slightly out-of-date hardware technology (0.35 μm) in order to minimize the production cost (less than five dollars) and increase security [34]. Furthermore, up to now, there was no real need for large memories in smartcard applications such as the holder's identification. According to major smartcard providers, the market pressure generated by emerging large storage demanding applications will lead to a rapid increase of the smartcard storage capacity. This evolution is however constrained by the smartcard tiny die size fixed to 25 mm² in the ISO standard [23], which pushes for more integration. This limited size is due to security considerations (to minimize the risk of physical attack [5]) and practical constraints (e.g., the chip should not break when the smartcard is flexed). Another solution to relax the storage limit is to extend the smartcard storage capacity with external memory modules. This is being done by Gemplus which recently announced Pinocchio [16], a smartcard equipped with 2 MB of Flash memory linked to the microcontroller by a bus. Since hardware security can no longer be provided on this memory, its content must be either non-sensitive or encrypted.

Another important issue is the performance of stable memory. Possible alternatives to the EEPROM are Flash memory and Ferroelectric RAM (FeRAM) [15] (see Table 2 for performance comparisons). Flash is more compact than EEPROM and represents a good candidate for high capacity smartcards [16]. However, Flash banks need to be erased before writing, which is extremely slow. This makes Flash memory appropriate for applications with a high read/write ratio (e.g., address books). FeRAM is undoubtedly an interesting option for smartcards as read and write times are both fast. Although its theoretical foundation was set in the early 1950s, FeRAM is just emerging as an industrial solution. Therefore, FeRAM is expensive, less secure than EEPROM or Flash, and its integration with traditional technologies (such as CPUs) remains an

Table 2. Performance of stable memories for the smartcard

Memory type	EEPROM	FLASH	FeRAM
Read time (/word)	60 to 150 ns	70 to 200 ns	150 to 200 ns
Write time (/word)	1 to 5 ms	5 to 10 μ s	150 to 200 ns
Erase time (/bank)	None	500 to 800 ms	None
Lifetime ^(*) (/cell)	10 ⁵ write cycles	10 ⁵ erase cycles	10 ¹⁰ to 10 ¹² write cycles

* A memory cell can be overwritten a finite number of time.

issue. Thus FeRAM could be considered a serious alternative only in the very long term [15].

Given these considerations, we assume in this paper a smartcard with a reasonable stable storage area (a few megabytes of EEPROM²) and a small RAM area (some kilobytes). Indeed, there is no clear interest in having a large RAM area, given that the smartcard is not autonomous, thus precluding asynchronous write operations. Moreover, more RAM means less EEPROM as the chip size is limited.

3.2 Impact on the PicoDBMS architecture

We now analyze the impact of the smartcard constraints on the PicoDBMS architecture, thus justifying why traditional database techniques, and even lightweight DBMS techniques, are irrelevant. The smartcard's properties and their impact are:

- *Highly secure*: smartcard's hardware security makes it the ideal storage support for private data. The PicoDBMS must contribute to the data security by providing access right management and a view mechanism that allows complex view definitions (i.e., supporting data composition and aggregation). The PicoDBMS code must not present security holes due to the use of sophisticated algorithms³.
- *Highly portable*: the smartcard is undoubtedly the most portable personal computer (the wallet computer). The data located on the smartcard are thus highly available. They are also highly vulnerable since the smartcard can be lost, stolen or accidentally destroyed. The main consequence is that durability cannot be enforced locally.
- *Limited storage resources*: despite the foreseen increase in storage capacity, the smartcard will remain the lightest representative of personal computers for a long time. This means that specific storage models and execution techniques must be devised to minimize the volume of persistent data (i.e., the database) and the memory consumption during execution. In addition, the functionalities of the PicoDBMS must be carefully selected and their implementation must be as light as possible. The lightest the PicoDBMS, the biggest the onboard database.
- *Stable storage is main memory*: smartcard stable memory provides the read speed and direct access granularity of a main memory. Thus, a PicoDBMS can be considered as a *main memory DBMS (MMDBMS)*. However the dramatic cost of writes distinguishes a PicoDBMS from a traditional MMDBMS. This impacts on the storage and access

² Considering Flash instead of EEPROM will not change our conclusions. It will just exacerbate them.

³ Most security holes are the results of software bugs [34].

methods of the PicoDBMS as well as the way transaction atomicity is achieved.

- *Non-autonomous*: compared to other computers, the smartcard has no independent power supply, thereby precluding disconnected and asynchronous processing. Thus, all transactions must be completed while the card is inserted in a terminal (unlike PDA, write operations cannot be cached in RAM and reported on stable storage asynchronously).

3.3 Problem statement

To summarize, our goal is to design a PicoDBMS including the following components:

- *Storage manager*: manages the storage of the database and the associated indices.
- *Query manager*: processes query plans composed of select, project, join, and aggregates.
- *Transaction manager*: enforces the ACID properties and participates to distributed transactions.
- *Access right manager*: provides access rights on base data and on complex user-defined views.

Thus, the PicoDBMS hosted in the chip provides the minimal subset of functionality that is strictly needed to manage in a secure way the data shared by all onboard applications. Other components (e.g., the GUI, a sort operator, etc.) can be hosted in the terminal or be dynamically downloaded when needed, without threatening security. In the rest of this paper, we concentrate on the components which require non-traditional techniques (storage manager, query manager, and transaction manager) and ignore the access right manager for which traditional techniques can be used.

When designing the PicoDBMS's components, we must follow several design rules derived from the smartcard's properties:

- *Compactness rule*: minimize the size of data structures and the PicoDBMS code to cope with the limited stable memory area (a few megabytes).
- *RAM rule*: minimize the RAM usage given its extremely limited size (some kilobytes).
- *Write rule*: minimize write operations given their dramatic cost (≈ 1 ms/word).
- *Read rule*: take advantage of the fast read operations (≈ 100 ns/word).
- *Access rule*: take advantage of the low granularity and direct access capability of the stable memory for both read and write operations.
- *Security rule*: never externalize private data from the chip and minimize the algorithms' complexity to avoid security holes.

4 PicoDBMS storage model

In this section, following the design rules for a PicoDBMS, we discuss the storage issues and propose a very compact model based on a combination of flat storage, domain storage, and ring storage. We also evaluate the storage cost of our storage model.

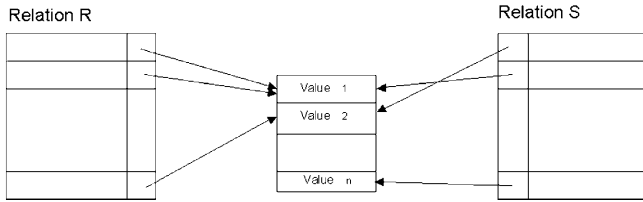


Fig. 1. Domain storage

4.1 Flat storage

The simplest way to organize data is *flat storage (FS)*, where tuples are stored sequentially and attribute values are embedded in the tuples. Although it does not impose it, the SCQL standard [24] considers FS as the reference storage model for smartcards. The main advantage of FS is access locality. However, in our context, FS has two main drawbacks:

- *Space consuming*: while normalization rules preclude attributes conjunction redundancy to occur, they do not avoid attribute value duplicates (e.g., the attribute *Doctor.Specialty* may contain many duplicates).
- *Inefficient*: in the absence of index structures, all operations are computed sequentially. While this is convenient for old fashion cards (some kilobytes of storage and a mono-relation select operator), this is no longer acceptable for future cards where storage capacity is likely to exceed 1 MB and queries can be rather complex.

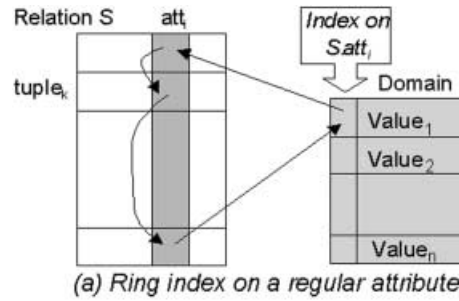
Adding index structures to FS may solve the second problem while worsening the first one. Thus, FS alone is not appropriate for a PicoDBMS.

4.2 Domain storage

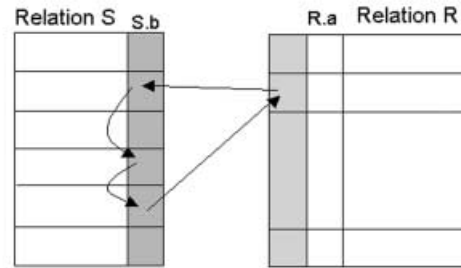
Based on the critique of FS, it follows that a PicoDBMS storage model should guarantee both data and index compactness. Let us first deal with data compactness. Since locality is no longer an issue in our context, pointer-based storage models inspired by MMDBMS [3, 27, 31] can help reducing the data storage cost. The basic idea is to preclude any duplicate value from occurring. This can be achieved by grouping values in domains (sets of unique values). We call this model *domain storage (DS)*. As shown in Fig. 1, tuples reference their attribute values by means of pointers. Furthermore, a domain can be shared among several attributes. This is particularly efficient for enumerated types, which vary on a small and determined set of values⁴.

One may wonder about the cost of tuple creation, update, and deletion since they may generate insertion and deletion of values in domains. While these actions are more complex than their FS counterpart, their implementation remains more efficient in the smartcard context, simply because the amount of data to be written is much smaller. To amortize the slight overhead of domain storage, we only store by domain all large attributes (i.e., greater than a pointer size) containing duplicates. Obviously, attributes with no duplicates (e.g., keys) need

⁴ Compression techniques can be advantageously used in conjunction with DS to increase compactness [17].



(a) Ring index on a regular attribute



(b) Ring index on a foreign key attribute

Fig. 2. Ring storage

not be stored by domain but with FS. Variable-size attributes – generally larger than a pointer – can also be advantageously stored in domains even if they do not contain duplicates. The benefit is not storage savings but memory management simplicity (all tuples of all relations become fixed-size) and log compactness (see Sect. 6).

4.3 Ring storage

We now address index compactness along with data compactness. Unlike disk-based DBMS that favor indices which preserve access locality, smartcards should make intensive use of secondary (i.e., pointer-based) indices. The issue here is to make these indices as compact as possible. Let us first consider select indices. A select index is typically made of two parts: a collection of values and a collection of pointers linking each value to all tuples sharing it. Assuming the indexed attribute varies on a domain, the index’s collection of values can be saved since it exactly corresponds to the domain extension. The extra cost incurred by the index is then reduced to the pointers linking index values to tuples.

Let us go one step further and get these pointers almost for free. The idea is to store these *value-to-tuple* pointers in place of the *tuple-to-value* pointers within the tuples (i.e., pointers stored in the tuples to reference their attribute values in the domains). This yields to an index structure which makes a ring from the domain values to the tuples. Hence, we call it *ring index* (see Fig. 2a). However, the ring index can also be used to access the domain values from the tuples and thus serve as data storage model. Thus we call *ring storage (RS)* the storage of a domain-based attribute indexed by a ring. The index storage cost is reduced to its lowest bound, that is, one pointer per domain value, whatever the cardinality of the indexed relation. This important storage saving is obtained at the price of extra work for projecting a tuple to the corresponding attribute since retrieving the value of a ring stored attribute means traversing

on average half of the ring (i.e., up to reaching the domain value).

Join indices [40] can be treated in a similar way. A join predicate of the form $(R.a = S.b)$ assumes that $R.a$ and $S.b$ vary on the same domain. Storing both $R.a$ and $S.b$ by means of rings leads to defining a join index. In this way, each domain value is linked by two separate rings to all tuples from R and S sharing the same join attribute value. However, most joins are performed on key attributes, $R.a$ being a primary key and $S.b$ being the foreign key referencing $R.a$. In our model, key attributes are not stored by domain but with FS. Nevertheless, since $R.a$ is the primary key of R , its extension forms precisely a domain, even if not stored outside of R . Since attributes $S.b$ take their values in $R.a$'s domain, they reference $R.a$ values by means of pointers. Thus, the domain-based storage model naturally implements for free a *unidirectional join index* from $S.b$ to $R.a$ (i.e., each S tuple is linked by a pointer to each R tuple matching with it). If traversals from $R.a$ to $S.b$ need to be optimized too, a *bi-directional join index* is required. This can be simply achieved by defining a ring index on $S.b$. Figure 2b shows the resulting situation where each R tuple is linked by a ring to all S tuples matching with it and vice versa. The cost of a bi-directional join index is restricted to a single pointer per R tuple, whatever the cardinality of S . Note that this situation resembles the well-known Codasyl model.

4.4 Storage cost evaluation

Our storage model combines FS, DS, and RS. Thus, the issue is to determine the best storage for each attribute. If the attributes need not be indexed, the choice is obviously between FS and DS. Otherwise, the choice is between RS and FS with a traditional index. Thus, we compare the storage cost for a single attribute, indexed or not, for each alternative. We introduce the following parameters:

- *CardRel*: cardinality of the relation holding the attribute.
- *a*: average length of the attribute (expressed in bytes).
- *p*: pointer size (3 bytes will be required to address “large” memory of future cards).
- *S*: selectivity factor of the attribute. $S = CardDom / CardRel$, where *CardDom* is the cardinality of the attribute domain extension (in all models). *S* measures the redundancy of the attribute (i.e., the same attribute value appears in $1/S$ tuples).

$$Cost(FS) = CardRel * a \quad // \text{ attribute storage cost in } // \text{ the relation}$$

$$Cost(DS) = CardRel * p \quad // \text{ attribute storage cost in } // \text{ the relation}$$

$$+ S * CardRel * a \quad // \text{ values storage cost in } // \text{ the domain}$$

$$Cost(Indexed.FS) = Cost(FS) \quad // \text{ flat attribute storage cost}$$

$$+ S * CardRel * a \quad // \text{ value storage cost in the } // \text{ index}$$

$$+ CardRel * p \quad // \text{ pointer storage cost in } // \text{ the index}$$

$$Cost(RS) = Cost(DS) \quad // \text{ domain-based attribute } // \text{ storage cost}$$

$$+ S * CardRel * p \quad // \text{ pointer storage cost in } // \text{ the index}$$

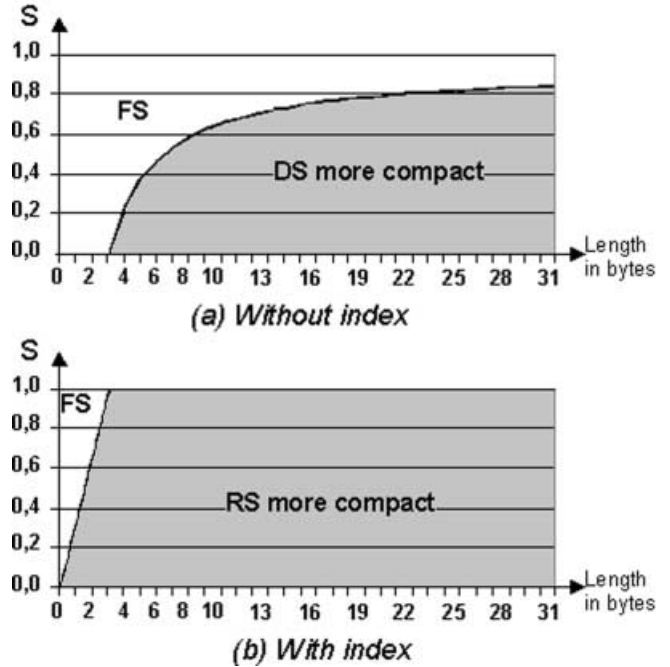


Fig. 3. Storage models' tradeoff

The cost equality between FS and DS gives: $S = (a-p)/a$.
 The cost equality between Indexed.FS and RS gives:

$$S = a/p$$

Figure 3a shows the different values of S and a for which FS and DS are equivalent. Thus, each curve divides the plan into a gain area for FS (above the curve) and a gain area for DS (under the curve). For values of a less than 3 (i.e., the size of a pointer), FS is obviously always more compact than DS. For higher values of a , DS becomes rapidly more compact than FS except for high values of S . For instance, considering $S = 0.5$, that is the same value is shared by only two tuples, DS outperforms FS for all a larger than 6 bytes. The higher a and the lower S , the better DS. The benefit of DS is thus particularly important for enumerated type attributes. Figure 3b compares Indexed.FS with RS. The superiority of RS is obvious, except for 1- and 2-byte-long key attributes. Thus, Figs. 3a and 3b are guidelines for the database designer to decide how to store each attribute, by considering its size and selectivity.

5 Query processing

Traditional query processing strives to exploit large main memory for storing temporary data structures (e.g., hash tables) and intermediate results. When main memory is not large enough to hold some data, state-of-the-art algorithms (e.g., hybrid hash join [33]) resort to materialization on disk to avoid memory overflow. These algorithms cannot be used for a PicoDBMS because:

- Given the write rule and the lifetime of stable memory, writes in stable memory are proscribed, even for temporary materialization.

- Dedicating a specific RAM area does not help since we cannot estimate its size a priori. Making it small increases the risk of memory overflow, thereby leading to writes in stable memory. Making it large reduces the stable memory area, already limited in a smartcard (RAM rule). Moreover, even a large RAM area cannot guarantee that query execution will not produce memory overflow [9].
- State-of-the-art algorithms are quite sophisticated, which precludes their implementation in a PicoDBMS whose code must be simple, compact, and secure (compactness and security rules).

To solve this problem, we propose query processing techniques that do not use any working RAM area nor incur any writes in stable memory. In the following, we describe these techniques for simple and complex queries, including aggregation and remove duplicates. We show the effectiveness of our solution through a performance analysis.

5.1 Basic query execution without RAM

We consider the execution of *SPJ* (*Select/Project/Join*) queries. Query processing is classically done in two steps. The query optimizer first generates an “optimal” *query execution plan* (*QEP*). The QEP is then executed by the query engine which implements an *execution model* and uses a library of relational operators [17]. The optimizer can consider different shapes of QEP: *left-deep*, *right-deep* or *bushy trees* (see Fig. 4). In a left-deep tree, operators are executed sequentially and each intermediate result is materialized. On the contrary, right-deep trees execute operators in a pipeline fashion, thus avoiding intermediate result materialization. However, they require materializing in memory all left relations. Bushy trees offer opportunities to deal with the size of intermediate results and memory consumption [38].

In a PicoDBMS, the query optimizer should not consider any of these execution trees as they incur materialization. The solution is to only use pipelining with *extreme right-deep trees* where all the operators (including select) are pipelined. As left operands are always base relations, they are already materialized in stable memory, thus allowing us to execute a plan with no RAM consumption. Pipeline execution can be easily achieved using the well-known *Iterator Model* [17]. In this model, each operator is an *iterator* that supports three procedure calls: *open* to prepare an operator for producing an item, *next* to produce an item, and *close* to perform final clean-up. A *QEP* is activated starting at the root of the operator tree and progressing towards the leaves. The dataflow in the model is demand-driven: a child operator passes a tuple to its parent node in response to a *next* call from the parent.

Let us now detail how select, project, and join are performed. These operators can be executed either sequentially or with a ring index. Given the access rule, the use of indices seems always to be the right choice. However, extreme right-deep trees allow us to speed-up a single select on the first base relation (e.g., *Drug.type* in our example), but using a ring index on the other selected attributes (e.g., *Visit.date*) may slow down execution as the rings need to be traversed to retrieve their value. Project operators are pushed up to the tree since no materialization occurs. Note that the final project incurs

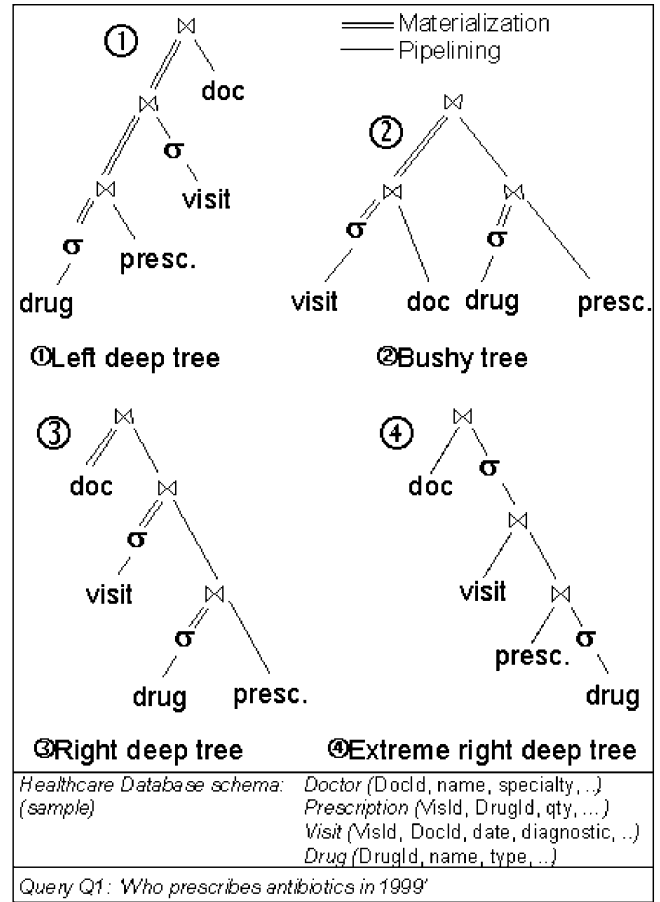


Fig. 4. Several execution trees for query Q1

an additional cost in case of ring attributes. Without indices, joining relations is done by a nested-loop algorithm since no other join technique can be applied without ad hoc structures (e.g., hash tables) and/or working area (e.g., sorting). The cost of indexed joins depends on the way indices are traversed. Consider the indexed join between *Doctor* (*n* tuples) and *Visit* (*m* tuples) on their key attribute. Assuming a unidirectional index, the join cost is proportional to $n * m$ starting with *Doctor* and to m starting with *Visit*. Assuming now a bi-directional index, the join cost becomes proportional to $n + m$ starting with *Doctor* and to $m^2 / 2n$ starting with *Visit* (retrieving the doctor associated to each visit incurs traversing half of a ring in average). In the latter case, a naïve nested loop join can be more efficient if the ring cardinality is greater than the target relation cardinality (i.e., when $m > n^2$). In that case, the database designer must clearly choose a unidirectional index between the two relations.

5.2 Complex query execution without RAM

We now consider the execution of aggregate, sort, and duplicate removal operators. At first glance, pipeline execution is not compatible with these operators which are classically performed on materialized intermediate results. Such materialization cannot occur either in the smartcard due to the RAM rule or in the terminal due to the security rule. Note that sorting can be done in the terminal since the output order of the

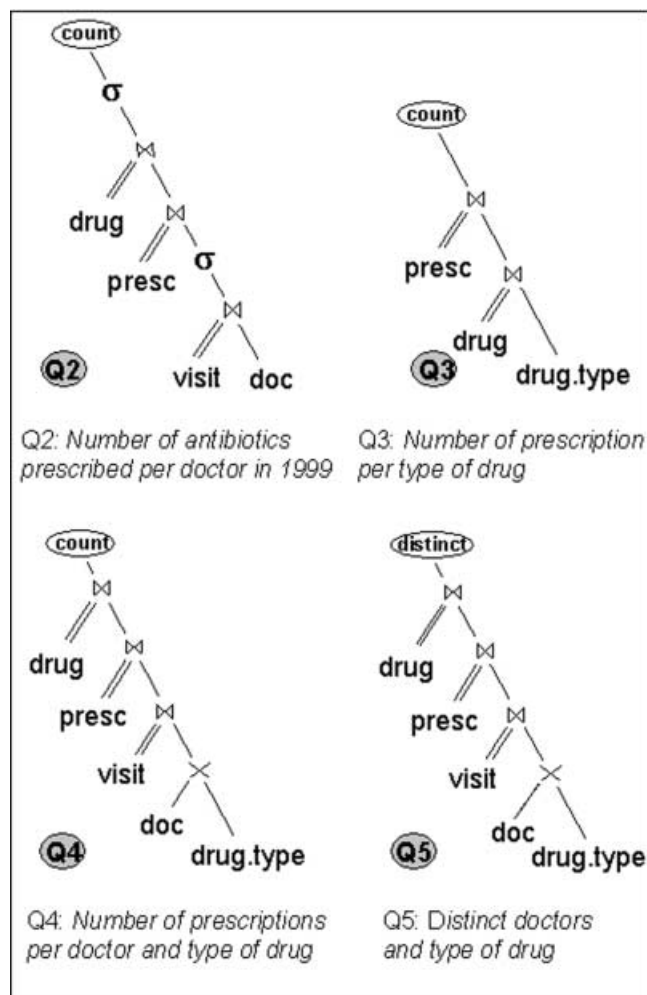


Fig. 5. Four ‘complex’ query execution plans

result tuples is not significant, i.e., depends on the DBMS algorithms.

We propose a solution to the above problem by exploiting two properties: (i) aggregate and duplicate removal can be done in pipeline if the incoming tuples are still grouped by distinct values; and (ii) pipeline operators are order-preserving since they consume (and produce) tuples in the arrival order. Thus, enforcing an adequate consumption order at the leaf of the execution tree allows pipelined aggregation and duplicate removal. For instance, the extreme right-deep tree of Fig. 4 delivers the tuples naturally grouped by *Drug.id*, thus allowing group queries on that attribute.

Let us now consider query Q2 of Fig. 5. As pictured, executing Q2 in pipeline requires rearranging the execution tree so that relation *Doctor* is explored first. Since *Doctor* contains distinct doctors, the tuples arriving to the *count* operator are naturally grouped by doctors.

The case of Q3 is harder. As the data must be grouped by *type of drugs* rather than by *Drug.id*, an additional join is required between relation *Drug* and domain *drug.type*. Domain values being unique, this join produces the tuples in the adequate order. If domain *Drug.type* does not exist, an operator must be introduced to sort relation *Drug* in pipeline. This can be done by performing n passes on *Drug* where n is the number of distinct values of *Drug.type*.

The case of Q4 is even trickier. The result must be grouped on two attributes (*Doctor.id* and *Drug.type*), introducing the need to start the tree with both relations! The solution is to insert a Cartesian product operator at the leaf of the tree in order to produce tuples ordered by *Doctor.id* and *Drug.type*. In this particular case, the query response time should be approximately n times greater than the same query without the ‘group by’ clause, where n is the number of distinct *types of drugs*.

Q5 retrieves the distinct couples of *doctor* and *type of prescribed drugs*. This query can be made similar to Q4 by expressing the distinct clause as an aggregate without function (i.e., the query “*select distinct a₁, . . . , a_n from . . .*” is equivalent to “*select a₁, . . . , a_n from . . . group by a₁, . . . , a_n*”). The unique difference is that the computation for a given group, i.e., (*distinct result tuple*) can stop as soon as one tuple has been produced.

5.3 Query optimization

Heuristic optimization is attractive. However, well-known heuristics such as processing select and project first do not work here. Using extreme right-deep trees makes the former impractical and invalidates the latter. Heuristics for join ordering are even more risky considering our data structures. Conversely, there are many arguments for an exhaustive search of the best plan. First, the search space is limited since: (i) there is a single algorithm for each operator, depending on the existing indices; (ii) only extreme right-deep trees are considered; and (iii) typical queries will not involve many relations. Second, exhaustive search using depth-first algorithms do not consume any RAM. Finally, exhaustive algorithms are simple and compact (even if they iterate a lot). Under the assumption that query optimization is required in a PicoDBMS, the remarks above strongly argue in favor of an exhaustive search strategy.

5.4 Performance evaluation

Our proposed query engine can handle fairly complex queries, taking advantage of the read and access rules⁵ while satisfying the compactness, write, RAM, and security rules. We now evaluate whether the PicoDBMS performance matches the smartcard application’s requirements, that is, any query issued by the application can be performed in reasonable time (i.e., may not exceed the user’s patience). Since the PicoDBMS code’s simplicity is an important consideration to conform to the compactness and security rules, we must also evaluate which acceleration techniques (i.e., ring indices, query optimization) are really mandatory. For instance, an accelerator reducing the response time from 10 ms to 1 ms is useless in the smartcard context⁶. Thus, unlike traditional performance evaluation, our major concern is on absolute rather than relative performance.

⁵ With traditional DBMS, such techniques will induce so many disk accesses that the system would thrash!

⁶ With traditional DBMS, such acceleration can improve the transactional throughput.

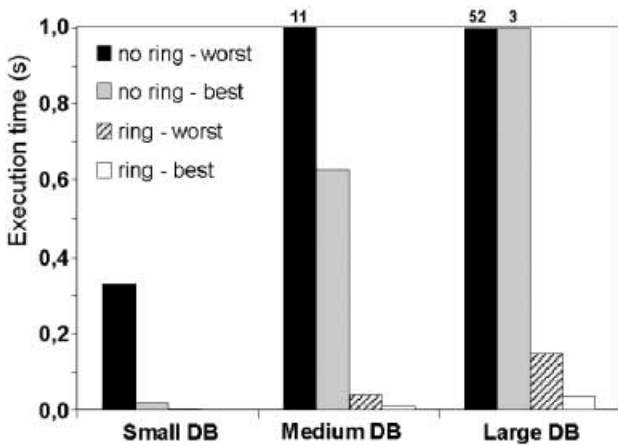


Fig. 6. Performance results for Q1

Evaluating absolute response time is complex in the smartcard environment because all platform parameters (e.g., processor speed, caching strategy, RAM, and EEPROM speed) strongly impact on the measurements⁷. Measuring the performance of our PicoDBMS on Bull’s smartcard technology is attractive but introduces two problems. First, Bull’s smartcards compatible with database applications are still prototypes [39]. Second, we are interested in providing the most general conclusions (i.e., as independent as possible of smartcard architectures). Therefore, we prefer to measure our query engine on two oldfashioned computers (a PC 486/25 Mhz and a Sun SparcStation 1+) which we felt roughly similar to forthcoming smartcard architectures. For each computer, we vary the system parameters (clock frequency, cache) and perform the experimentation tests. The performance ratios between all configurations were roughly constant (i.e., whatever the query), the slowest configuration (Intel 486 with no cache) performing eight times worse than the fastest (RISC with cache). In the following, we present response times for the slowest architecture to check the viability of our solutions in the worst environment.

We generated three instances of a simplified healthcare database: the *small*, *medium*, and *large* databases containing, respectively, (10, 30, 50) doctors, (100, 500, 1,000) visits, (300, 2,000, 5,000) prescriptions, and (40, 120, 200) drugs. Although we tested several queries, we describe below only the two most significant. Query Q1, which contains three joins and two selects on *Visit* and *Drug* (with selectivities of 20% and 5%), is representative of medium-complexity queries. Query Q4, which performs an aggregate on two attributes and requires the introduction of a Cartesian product, is representative of complex queries. For each query, we measure the performance for all possible query execution plans, excluding those which induce additional Cartesian product, varying the storage choices (with and without select and join ring indices). Figures 6 and 7 show the results for both best and worst plans on databases built with or without join indices.

Considering SPJ queries, the PicoDBMS performance clearly matches the application’s requirements as soon as join rings are used. Indeed, the performance with join rings is at

⁷ With traditional DBMS, very slow disk access allows us to ignore finer parameters.

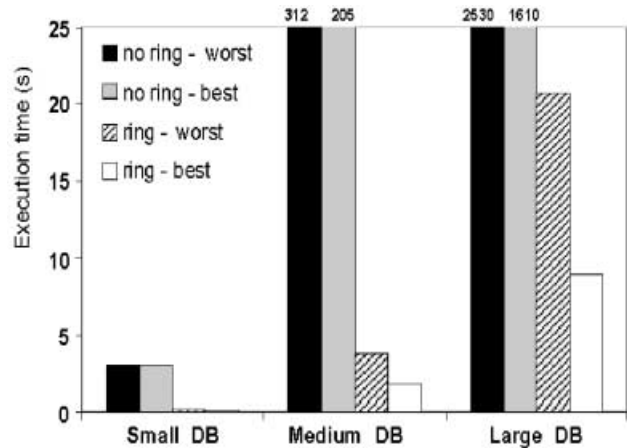


Fig. 7. Performance results for Q4

most 146 ms for the largest database and with the worst execution plan. With small databases, all the acceleration techniques can be discarded, while with larger ones, join rings remain necessary to obtain good response time. In that case, the absolute gain (110 ms) between the best and the worst plan does not justify the use of a query optimizer.

The performance of aggregate queries is clearly the worst because they introduce a Cartesian product at the leaf of the execution tree. Join rings are useful for medium and large databases. With large databases, the optimizer turns out to be necessary since the worst execution plan with join rings achieves a rather long response time (20.6 s).

The influence of ring indices for selects (not shown) is insignificant. Depending on the selectivity, it can bring slight improvement or overhead on the results. Although it may achieve an important relative speed-up for the select itself, the absolute gain is not significant considering the small influence of select on the global query execution cost (which is not the case in disk-based DBMS). Select ring indices are, however, useful for queries with aggregates or duplicate removal, that can result in a join between a relation and the domain attribute. In that case, the select index plays the role of a join index, thereby generating a significant gain on large relations and large domains.

Thus, this performance evaluation shows that our approach is feasible and that join indices are mandatory in all cases while query optimization turns out to be useful only with large databases and complex queries.

6 Transaction management

Like any data server, a PicoDBMS must enforce the well-known transactional ACID properties [8] to guarantee the consistency of the local data it manages as well as be able to participate in distributed transactions. We discuss below these properties with respect to a PicoDBMS.

- *Atomicity: local atomicity* means that the set of actions performed by the PicoDBMS on a transaction’s behalf is made persistent following the *all or nothing* scheme. *Global atomicity*: this means that all data servers – including the PicoDBMS – accessed by a distributed transaction

agree on the same transaction outcome (either commit or rollback). The distinguishing features of a PicoDBMS regarding atomicity are no demarcation between main memory and persistent storage, the dramatic cost of writes, and the fact that they cannot be deferred.

- *Consistency*: this property ensures that the actions performed by the PicoDBMS satisfy all integrity constraints defined on the local data. Considering that traditional integrity constraint management can be used, we do not discuss it any further.
- *Isolation*: this property guarantees the serializability of concurrent executions. A PicoDBMS manages personal data and is typically single-user⁸. Furthermore, smartcard operating systems do not even support multithreading. Therefore, isolation is useless here.
- *Durability*: durability means that committed updates are never lost whatever the situation (i.e., even in case of a media failure). Durability cannot be enforced locally by the PicoDBMS because the smartcard is more likely to be stolen, lost or destroyed than a traditional computer. Indeed, mobility and smallness play against safety. Consequently, durability must be enforced through the network. The major issue is then preserving the privacy of data while delegating the durability to an external agent.

The remainder of this section addresses local atomicity, global atomicity, and durability.

6.1 Local atomicity

There are basically two ways to perform updates in a DBMS. The updates are either performed on *shadow objects* that are atomically integrated in the database at commit time or done *in place* (i.e., the transaction updates the shared copy of the database objects) [8]. We discuss these two traditional models below.

- *Shadow update*: This model is rarely employed in disk-based DBMSs because it destroys data locality on disk and increases concurrent updates on the catalog. In a PicoDBMS, disk locality and concurrency are not a concern. This model has been shown to be convenient for smartcards equipped with a small Flash memory [25]. However, it is poorly adapted to pointer-based storage models like RS since the object location changes at every update. In addition, the cost incurred by shadowing grows with the memory size. Indeed, either the granularity of the shadow objects increases or the paths to be duplicated in the catalog become longer. In both cases, the writing cost – which is the dominant factor – increases.
- *Update in-place*: write-ahead logging (WAL) [8] is required in this model to undo the effects of an aborted transaction. Unfortunately, the relative cost of WAL is much higher in a PicoDBMS than in a traditional disk-based DBMS which uses buffering to minimize I/Os. In a smartcard, the log must be written for each update since each update becomes immediately persistent. This roughly doubles the cost of writing.

⁸ Even if the data managed by the PicoDBMS are shared among multiple users (e.g., as in the healthcare application), the PicoDBMS serves a single user at a time.

Despite its drawbacks, *update in-place* is better suited than *shadow update* for a PicoDBMS because it accommodates pointer-based storage models and its cost is insensitive to the rapid growth of stable memory capacity. We also propose two optimizations to *update in-place*:

- *Pointer-based logging*: traditional WAL logs the values of all modified data. RS allows a finer granularity by logging pointers in place of values. The smallest the log records, the cheapest the WAL. The logging process must consider two types of information:
 - *Values*: in case of a tuple update, the log record must contain the tuple address and the old attribute values, that is a pointer for all RS stored attributes and a regular value for FS stored attributes. In case of a tuple insertion or deletion, assuming each tuple header contains a status bit (i.e., dead or alive), only the tuple address has to be logged in order to recover its state.
 - *Rings*: tuple insertion, deletion, and update (of a ring attribute) modify the structure of each ring traversing the corresponding tuple t . Since a ring is a circular chain of pointers, recovering its state means recovering the *next* pointer of t 's predecessor (let us call it t_{pred}). The information to restore in $t_{pred.next}$ is either t 's address if t has been updated or deleted, or $t.next$ if t has been inserted. t 's address already belongs to the log (see above) and $t.next$ does not have to be logged since t 's content still exists in stable storage at recovery time. The issue is how to identify t_{pred} at recovery time. Logging this information can be saved at the price of traversing the whole ring starting from t , until reaching t again. Thus, ring recovery comes for free in terms of logging.
- *Garbage-collecting values*: insertion and deletion of domain values (domain values are never modified) should be logged as any other updates. This overhead can be avoided by implementing a deferred garbage collector that destroys all domain values no longer referenced by any tuple. Garbage-collecting a domain amounts to execute an ad hoc semi-join operator between the domain and all relations varying on it which discards the domain values that do not match⁹. The benefit of this solution is threefold: (i) the lazy deletion of unreferenced values does not entail the storage model coherency; (ii) garbage-collecting domain values is required anyway by RS (even in the absence of transaction control); and (iii) a deferred garbage-collector can be implemented without reference counters, thereby saving storage space. The deferred garbage collector cannot work in the background since smartcards do not yet support multi-threading. The most pragmatic solution is to launch it manually when the card is nearly full. An alternative to this manual procedure is to execute the garbage collector automatically at each card connection on a very small subset of the database (so that its cost remains hidden to the user). Garbage-collecting the database in such an incremental way is straightforward since domain values are examined one after the other.

⁹ Unlike reachability algorithms that start from the persistent roots and need marking [6], the proposed garbage-collector starts from the persistent leaves (i.e., the domain values) and exploits them one after the other, in a pipelined fashion (thus, it conforms to the RAM rule).

The update in-place model along with pointer-based logging and deferred garbage-collector reduces logging cost to its lowest bound, that is, a tuple address for inserted and deleted tuples, and the values of updated attributes (again, a pointer for DS and RS stored attributes).

6.2 Global atomicity

Global atomicity is traditionally enforced by an *atomic commitment protocol (ACP)*. The most well known and widely used ACP is 2PC [8]. While extensively studied [19] and standardized [21, 29, 41], 2PC suffers from the following weaknesses in our context:

- *Need for a standard prepared state*: any server must externalize the standard *Xa* interface [41] to participate to 2PC. Unfortunately, ISO defines a transactional interface for smartcards but it does not cover distributed transactions [24]. In addition, participating to 2PC requires building a local prepared state that consumes valuable resources.
- *Disconnection means aborting*: a smartcard can be extracted from its terminal or its mobile host (e.g., a cellular phone) can be temporarily unreachable during 2PC. A participant's disconnection leads 2PC to abort the transaction even if all its operations have been successfully executed.
- *Badly adapted to moving participants*: the 2PC incurs two message rounds to commit a transaction. Considering the high cost of wireless communication, the overhead is significant for mobile terminals equipped with a smartcard reader (e.g., PDA, cellular phones).

As its name indicates, 2PC has two phases: the *voting* phase and the *decision* phase. The voting phase is the means by which the coordinator checks whether or not the participants can locally guarantee the ACID properties of the distributed transaction. The decision is *commit* if all participants vote *yes* and *abort* otherwise. Thus, the voting phase introduces an uncertainty period at transaction termination that leads to the aforementioned drawbacks.

Variations of *one-phase commit* protocols (*1PC*) have been recently proposed [2, 4, 35]. As stated in [2], 1PC eliminates the voting phase of 2PC by enforcing the following properties on the participant's behavior: (1) all operations are acknowledged before the 1PC is launched; (2) there are no deferred integrity constraints; (3) all participants are ruled by a rigorous concurrency control scheduler; and (4) all updates are logged on stable storage before 1PC is launched. These assumptions guarantee, respectively, the A, C, I, D properties before the ACP is launched. Then, the ACP reduces to a single phase, that is broadcasting the coordinator's decision to all participants (this decision is commit if all transaction's operations have been successfully executed and abort otherwise). If a crash or a disconnection precludes a participant from conforming to this decision, the corresponding transaction branch is simply forward recovered (potentially at the next reconnection). While the assumptions on the participant's behavior seem constraining in the general case, they are quite acceptable in the smartcard context [10]. Property (1) is common to all ACPs and is enforced by the ISO7816 standard [22]; property (2) conforms to the fact that PicoDBMS have lighter capabilities

than full-fledged DBMS; and property (3) is satisfied by definition since smartcards do not support parallel executions. Property (4) is discussed in Sect. 6.3.

Eliminating the voting phase of the ACP solves altogether the three aforementioned problems. However, one may wonder about the interoperability between transaction managers and data managers supporting different protocols (either 1PC or 2PC). We have shown in [1] that the participation of legacy (i.e., 2PC compliant) data managers in 1PC is straightforward. Conversely, the participation of 1PC compliant data managers (e.g., a smartcard) in the 2PC can be achieved by associating a *log agent* to each participant. The role of the log agent is twofold. First, it manages the data manager's part of the 1PC's coordinator log, forces it to stable storage during the 2PC prepare phase, and exploits it if the transaction branch needs to be forward-recovered. Second, it translates the 2PC interface into that of 1PC. The log agent can be located on the terminal, so that the benefit of 1PC is lost for the terminal but it is preserved for the smartcard.

6.3 Durability

Most 1PC protocols assume that the coordinator is in charge of logging all participants' updates before triggering the ACP (all these protocols belong to the coordinator log family). *Coordinator log* [35] and *implicit yes vote* [4] assume that the participants piggyback their log records on the acknowledgment messages of each operation while *coordinator logical log* [2] assumes that the coordinator logs all operations sent to each participant. In all cases, the durability of the distributed transaction relies on the coordinator log. Thus, 1PC is a means by which global atomicity and durability can be solved altogether, at the same price.

Two issues remain to be solved: (i) where to store the coordinator log; and (ii) how to preserve the security rule, that is, how to make the log content as secure as the data stored in the smartcard. Since the log must sustain any kind of failure, it must be stored on the network by a trustee server (e.g., a public organism, a central bank, the card issuer, etc.). If some transactions are executed in disconnected mode (e.g., on a mobile terminal), the durability will be effective only at the time the terminal reconnects to the network. Protecting the log content against attacks imposes encryption. The way encryption is performed depends on the model of logging. If the coordinator log is fed by the log records piggybacked by the participants, the smartcard can encrypt them with an algorithm based on a private key (e.g., DES [28]). Otherwise (i.e., if the *coordinator logical log* scheme is selected), the smartcard can provide the coordinator with a public key that will be used by the coordinator itself to encrypt its log [32].

6.4 Transaction cost evaluation

The goal of this section is to approximate the time required by a representative update transaction. The objective is to confirm whether or not the write performance of smartcards assumed in this paper is acceptable for database applications like health cards. To this end, we estimate the time required to create a tuple in a relation, including the creation of domain values,

the insertion of the tuple in the rings potentially defined on this relation and the log time. Let us introduce the following parameters, in addition to those already defined in Sect. 4.4:

- $nbAttFS$: number of FS stored attributes
- $nbAttDS$: number of DS stored attributes
- $nbAttRS$: number of RS stored attributes
- w : size of a word (4 bytes in a 32-bit card)
- t : time to write one word in stable storage (5 ms in the worst case)

$$\begin{aligned} \text{Cost}(\text{insertTuple}) = & \\ & ((nbAttFS*a + nbAttDS*p + nbAttRS*p)/w) \quad // \textcircled{1} \\ & + (nbAttRS + nbAttDS) * S * [a/w] \quad // \textcircled{2} \\ & + nbAttRS * [p/w] \quad // \textcircled{3} \\ & + [p/w] \quad // \textcircled{4} \\ &) * t \quad // \textcircled{5} \end{aligned}$$

- ① Tuple size
- ② Domain values size. $S \approx$ probability to create a new domain value
- ③ Ring pointers to be updated
- ④ Log record size
- ⑤ Write time

Let us consider a representative transaction executed on the healthcare. This transaction inserts a new tuple in *Doctor* and *Visit* and five tuples in *Prescription* and *Drug*. This is somehow a worst case for this application in the sense that the visited doctor is a new one and prescribes five new drugs. The considered attribute distribution is as follows:

<i>Doctor</i>	($nbAttFS=3$, $nbAttDS=4$, $nbAttRS=0$),
<i>Visit</i>	($nbAttFS=2$, $nbAttDS=3$, $nbAttRS=2$),
<i>Prescription</i>	($nbAttFS=1$, $nbAttDS=1$, $nbAttRS=2$),
<i>Drug</i>	($nbAttFS=2$, $nbAttDS=4$, $nbAttRS=0$).

The average attribute length a is fixed to 10 bytes. Figure 8 plots the update transaction execution time depending on S ($S = 0$ means that all attribute values already exist in the domains, while $S = 1$ means that all these values need be inserted in the domains).

The figure is self-explanatory. Note that the logging cost represents less than 3% of the total cost. This simple analysis shows that the time expected for this kind of transaction (less than 1 s) is clearly compatible with the healthcare application's requirements.

7 Conclusion

As smartcards become more and more versatile, multi-application, and powerful, the need for database techniques arises. However, smartcards have severe hardware limitations which make traditional database technology irrelevant. The major problem is scaling down database techniques so they perform well under these limitations. In this paper, we addressed this problem and proposed the design of a PicoDBMS, concentrating on the components which require non-traditional techniques (storage manager, query manager, and transaction manager).

This paper makes several contributions. First, we analyzed the requirements for a PicoDBMS based on a healthcare application which is representative of personal

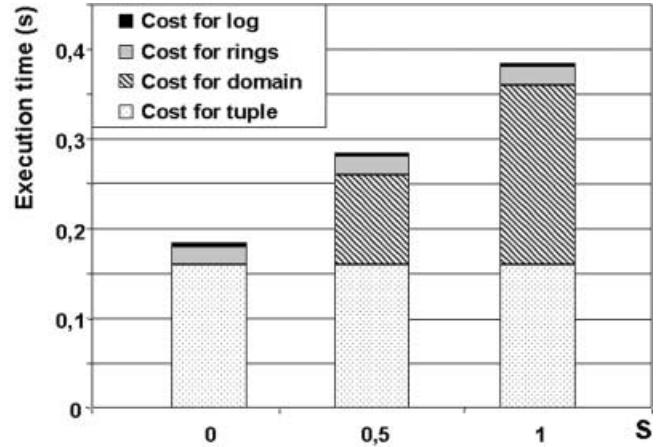


Fig. 8. Performance of a typical update transaction

folder applications and has strong database requirements. We showed that the minimal functionality should include select/project/join/aggregate, access right management, and views as well as transaction's atomicity and durability.

Second, we gave an in-depth analysis of the problem by considering the smartcard hardware trends. Based on this analysis, we assumed a smartcard with a reasonable stable memory of a few megabytes and a small RAM of some kilobytes, and we derived design rules for a PicoDBMS architecture.

Third, we proposed a new highly compact storage model that combines flat storage (FS), domain storage (DS), and ring storage (RS). Ring storage reduces the indexing cost to its lowest bound. Based on performance evaluation, we derived guidelines to decide the best way to store an attribute.

Fourth, we proposed query processing techniques which handle complex query plans with no RAM consumption. This is achieved by considering extreme right-deep trees which can pipeline all operators of the plan including aggregates. We also argued that, if query optimization is needed, the strategy should be exhaustive search. We measured the performance of our execution model with an implementation of our query engine on two old-fashioned computers which we configured to be similar to forthcoming smartcard architectures. We showed that the resulting performance matches the smartcard application's requirements.

Finally, we proposed techniques for transaction atomicity and durability. Local atomicity is achieved through update in-place with two optimizations which exploit the storage model: pointer-based logging and garbage collection of domain values. Global atomicity and durability are enforced by IPC which is easily applicable in the smartcard context and more efficient than 2PC. We showed that the performance of typical update transactions is acceptable for representative applications like the health card.

This work is done in the context of a new project with Bull Smart Cards and Terminals. The next step is to port our PicoDBMS prototype on Bull's smartcard new technology, called *OverSoft* [12], and to assess its functionality and performance on real-world applications. To this end, a benchmark dedicated to PicoDBMS must be set up. We also plan to address open issues such as protected logging for durability, query execution on encrypted data (e.g., stored in an external Flash), and statistics maintenance on a population of cards.

References

1. Abdallah M., Bobineau C., Guerraoui R., Pucheral P.: Specification of the transaction service. Esprit project OpenDREAMS-II n° 25262, Deliverable n° R13, 1998
2. Abdallah M., Guerraoui R., Pucheral P.: One-phase commit: does it make sense? Int. Conf. on Parallel and Distributed Systems (ICPADS), 1998
3. Ammann A., Hanrahan M., Krishnamurthy R.: design of a memory resident DBMS. IEEE COMPCON, 1985
4. Al-Houmailly Y., Chrysanthos P.K.: Two-phase commit in gigabit-networked distributed databases. Int. Conf. on Parallel and Distributed Computing Systems (PDCS), 1995
5. Anderson R., Kuhn M.: Tamper resistance – a cautionary note. USENIX Workshop on Electronic Commerce, 1996
6. Amsaleg L., Franklin M.J., Gruber O.: Efficient incremental garbage collection for client-server object database systems. Int. Conf. on Very Large Databases (VLDB), 1995
7. Bobineau C., Bouganim L., Pucheral P., Valduriez P.: PicoDBMS: scaling down database techniques for the smartcard (Best Paper Award). Int. Conf. on Very Large Databases (VLDB), 2000
8. Bernstein P.A., Hadzilacos V., Goodman N.: Concurrency control and recovery in database systems. Addison-Wesley, Reading, Mass., USA, 1987
9. Bouganim L., Kapitskaia O., Valduriez P.: Memory-adaptive scheduling for large query execution. Int. Conf. on Information and Knowledge Management (CIKM), 1998
10. Bobineau C., Pucheral P., Abdallah M.: A unilateral commit protocol for mobile and disconnected computing. Int. Conf. On Parallel and Distributed Computing Systems (PDCS), 2000
11. van Bommel F.A., Sembritzki J., Buettner H.-G.: Overview on healthcard projects and standards. Health Cards Int. Conf., 1999
12. Bull S.A.: Bull unveils iSimplify! the personal portable portal. Available at: http://www.bull.com:80/bull_news/
13. Carrasco L.C.: RDBMS's for Java cards? What a senseless idea! Available at: www.sqlmachine.com, 1999
14. DataQuest.: Chip card market and technology charge ahead. MSAM-WW-DP-9808, 1998
15. Dipert B.: FRAM: Ready to ditch niche? EDN Access Magazine, Cahners, London, 1997
16. Gemplus.: SIM Cards: From kilobytes to megabytes. Available at: www.gemplus.fr/about/pressroom/, 1999
17. Graefe G.: Query evaluation techniques for large databases. ACM Comput Surv, 25(2), 1993
18. Graefe G.: The new database imperatives. Int. Conf. on Data Engineering (ICDE), 1998
19. Gray J., Reuter A.: Transaction processing. Concepts and Techniques. Morgan Kaufmann, San Francisco, 1993
20. IBM Corporation.: DB2 Everywhere – administration and application programming guide. IBM Software Documentation, SC26-9675-00, 1999
21. International Standardization Organization (ISO).: Information technology - open systems interconnection - distributed transaction processing. ISO/IEC 10026, 1992
22. International Standardization Organization (ISO).: Integrated circuit(s) cards with contacts – part 3: electronic signal and transmission protocols. ISO/IEC 7816-3, 1997
23. International Standardization Organization (ISO).: Integrated circuit(s) cards with contacts – part 1: physical characteristics. ISO/IEC 7816-1, 1998
24. International Standardization Organization (ISO).: Integrated circuit(s) cards with contacts – part 7: interindustry commands for structured card query language (SCQL). ISO/IEC 7816-7, 1999
25. Lecomte S., Trane P.: Failure recovery using action log for smartcards transaction based system. IEEE Online Testing Workshop, 1997
26. Microsoft Corporation.: Windows for smartcards toolkit for visual basic 6.0. Available at: www.microsoft.com/windowsce/smartcard/, 2000
27. Missikov M., Scholl M.: Relational queries in a domain based DBMS. ACM SIGMOD Int. Conf. on Management of Data, 1983
28. National Institute of Standards and Technology.: Announcing the Data Encryption Standard (DES). FIPS PUB 46-2, 1993
29. Object Management Group.: Object transaction service. Document 94.8.4, OMG editor, 1994
30. Oracle Corporation.: Oracle 8i Lite - Oracle Lite SQL reference. Oracle documentation, A73270-01, 1999
31. Pucheral P., Thévenin J.M., Valduriez P.: Efficient main memory data management using the DBGraph storage model. Int. Conf. on Very Large Databases (VLDB), 1990
32. RSA Laboratories.: PKCS # 1: RSA Encryption Standard. RSA Laboratories Technical Note, v.1.5, 1993
33. Schneider D., DeWitt D.: A performance evaluation of four parallel join algorithms in a shared-nothing multiprocessor environment. ACM-SIGMOD Int. Conf., 1989
34. Schneier B., Shostack A.: Breaking up is hard to do: modeling security threats for smart cards. USENIX Symposium on Smart Cards, 1999
35. Stamos J., Cristian F.: A low-cost atomic commit protocol. IEEE Symposium on Reliable Distributed Systems, 1990
36. Sun Microsystems.: JavaCard 2.1 application programming interface specification. JavaSoft documentation, 1999
37. Sybase Inc.: Sybase adaptive server anywhere reference. CT75KNA, 1999
38. Shekita E., Young H., Tan K.L.: Multi-join optimization for symmetric multiprocessors. Int. Conf. on Very Large Data Bases (VLDB), 1993
39. Tual J.-P.: MASSC: a generic architecture for multiapplication smart cards. IEEE Micro J, N° 0272-1739/99, 1999
40. Valduriez P.: Join indices. ACM Trans. Database Syst, 12(2), 1987
41. X/Open.: Distributed transaction processing: reference model. X/Open Guide, Version 3. G307., X/Open Company Limited, 1996