

Copyright © 1979, by the author(s).
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

RETROSPECTION ON A DATA BASE SYSTEM

by

M. Stonebraker

Memorandum No. UCB/ERL M79/4

18 January 1979

RETROSPECTION OF A DATA BASE SYSTEM

by

Michael Stonebraker

Memorandum No. UCB/ERL M79/4

18 January 1979

ELECTRONICS RESEARCH LABORATORY

College of Engineering
University of California, Berkeley
94720

REQUIEM FOR A DATA BASE SYSTEM

by

Michael Stonebraker

DEPARTMENT OF ELECTRICAL ENGINEERING AND COMPUTER SCIENCE

UNIVERSITY OF CALIFORNIA

BERKELEY, CA.

ABSTRACT

This paper describes the implementation history of the INGRES data base system. It focuses on mistakes that were made in progress rather than on eventual corrections. Some attention is also given to the role of structured design in a data base system implementation and to the problem of supporting non trivial users. Lastly, miscellaneous impressions of UNIX, the PDP-11 and data models are given.

I INTRODUCTION

This paper was written in response to several requests to know what really happened in the INGRES data base management system project [STON76a] and why. To the extent that it contains practical wisdom for other implementation projects, it serves its purpose. To the extent that it is self-righteous defense of the existing design, the author apolo-

gizes in advance.

It may be premature to write such a document, since INGRES has only been fully operational for three years and user experience is still somewhat limited. Hence, the ultimate jury, real users, has not yet made a full report. The reason for reporting now is that we have reached a turning point. Until now, the goal was to make INGRES "really work", i.e. efficiently, reliably and without surprises (bugs) for users. There are now only marginal returns to pursuing that goal. Consequently, the project is taking new directions, which are discussed below.

This paper is organized as follows. In Section II we trace the history of the project through its various phases and highlight the more significant events that took place. Then, in Section III, we discuss several lessons that we had to learn the hard way. Section IV takes a critical look at the current design of INGRES and discusses some of the mistakes. Next, Section V consists of an assortment of random comments. Lastly, Section VI outlines the future plans of the project.

II HISTORY

The project can be roughly decomposed into three periods:

- 1) the early times -- 3/73-6/74
- 2) the first implementation -- 6/74-9/75

3) make it really work -- 9/75-present

We discuss each period in turn.

2.1 The Early Times

The project began in 1973 when Eugene Wong and I agreed to read and discuss literature relating to relational data bases. From the beginning we were both enthusiastic about an implementation. It did not phase either one of us that we possessed no experience whatsoever in leading a non trivial implementation effort. In fact, neither of us had ever written a sizeable computer program.

Our first task was to find a suitable machine environment for an implementation. It became clear quickly that no machine which we had access to was appropriate for an interactive data base system. Through various subterfuges (mainly engineered by Eugene Wong and Pravin Varaiya) we obtained about \$90,000. for hardware. The liability that we obtained was a commitment to write a geo-data system for the Urban Economics Group led by Pravin Varaiya and Roland Artle.

Our major concerns in selcting hardware were in obtaining large (50 or 100 megabytes at the time) disks and a decent software environment. After studying the UNIX CACM paper [RITC75], I was convinced that we should use UNIX and buy whatever hardware we could afford to make it run. We placed

a hardware order in February of 1974 and had a system in September of the same year.

We decided to offer a seminar running from September 1973 to June 1974 in which a design would be pursued. Somewhat symbiotically the seminar split into two groups: One, led by Gene, would plan the language; the other, led by me, would plan the support system. The language group converged quickly on the basic tenets of QUEL for retrieve operations. As soon as UNIX was chosen, my group laid out the system catalogs (data dictionary) and the access method interface. It never occurred to us that anyone would seriously consider making the data dictionary separate from the data base system (as is a common practice today). An idea from the very start had been to have several implementations of the access method interface. Each would have the same calling conventions for simplicity and would function interchangeably. We were committed to the relational principle that users see nothing of the underlying storage structure. Hence, no provisions were made to allow a user to access a lower level of the system (as is done in some other implementations).

During the winter of 1974 the notion of tuple substitution was developed as a method for "solving" queries. This notion of decomposition strongly influenced the resulting design. For example, having a level in the system that corresponded to the "one-variable query processor" occurred because decomposition required it.

In summary, the salient features of INGRES at the time were:

- 1) QUEL retrieval was defined
- 2) an integrated data dictionary was proposed
- 3) multiple implementations of the access methods were suggested
- 4) a "pure" relational system was agreed on
- 5) decomposition was developed

This first period ended with the delivery of a PDP-11 in June 1974 which could be used on an interim basis for code development. Hence, we could begin implementing before our own machine arrived. The project was organized as a chief programmer team of four persons under the direction of Gerry Held. This same organizational structure remains today.

2.2 The First Implementation

We expected to exploit the natural parallelism which multiple UNIX processes allowed. Hence, decomposition would be a process to run in parallel with the one-variable query processor (OVQP). The utilities (e.g. to create relations, destroy them and modify their storage structure) would be several overlays but nobody was exactly sure where they would go. By this time we had decided to take protection seriously and that a data base administrator was an appropriate concept. He or she should own all the physical UNIX files, and the INGRES object code should execute in protected mode. Because the terminal monitor allowed the

user to directly edit files, we had to protect the rest of INGRES from it. Hence, it had to be a separate process. The notion of query modification for protection, integrity control and views was developed during this time. It would be implemented with the parser but no thought was given to the form of this module. During the summer of 1974 the process structure changed several times. Moreover, no one could coherently check any code because everyone needed the access methods as part of their code and they did not work yet.

About this time another version of QUEL was developed which included updates and more general aggregates. This version survives today except for the keyword syntax, which was changed in early 1975.

By the end of the summer we had some access method code, some routines to access the data dictionary (to create and destroy relations for example), and a terminal monitor, along with pieces of DECOMP and OVQP. In September, the department arranged to invite Ken Thompson (the creator of UNIX in conjunction with Dennis Ritchie) to Berkeley for a two week visit. Ken was instrumental in getting UNIX to run on the INGRES machine and introduced us to YACC as a parser generator.

In January of 1975 we invited Ted Codd to come to Berkeley in early March to see a demonstration of INGRES. The final two weeks before his visit everyone worked night and day so

that we would have something to show him. What we demonstrated was a very "buggy" system with the following characteristics:

1) the access methods "sort of" worked. Retrieves worked on all five implementations of the access methods (heap, hash, compressed hash, index and compressed index). However, only heaps could be updated without fear of disaster.

2) decomposition was implemented by brute force

3) a primitive data base load program existed but few other services

4) all the messy interprocess problems had been ignored. For example, there was no way to reset INGRES so it would stop executing the current command and be ready to do something new. Instead of being able to flush all the processes, we simply killed them.

5) There were many bugs. For example, boolean operators sometimes worked incorrectly. The average function applied to a relation with no tuples produced a weird response, etc.

At this point it became clear that the punctuation oriented syntax for QUEL was horrible and it was scrapped in favor of a keyword oriented approach. The designers of SEQUEL saw this important point sooner than we did. This was the last significant change to the language.

During this time the "B-tree" debate raged. The pros and cons of dynamic and static directories were argued. We wrote the paper "B-trees Reexamined" [HELD78] during this period and believed its contents. This is one of the mistakes discussed in Section IV.

Lastly, it became clear that we needed a coupling to a host language. Moreover, "C" was the only possible candidate, since it alone allowed interprocess communication; a fact essential for INGRES operation. As a result we began work on a preprocessor EQUOL [ALLM76], to allow convenient access to INGRES from "C".

The end of this initial implementation period occurred when we acquired a user. Through Ken Thompson, to whom a tape of an early system had been sent, and through a group at Bell Labs in Holmdel, Mr. Dan Gielan of New York Telephone Co. became interested in using our system. After a trial period using our machine, he obtained his own and set about tailoring INGRES to his environment and fixing its flaws (many bugs, bad performance, no concurrency control, no recovery, shakey physical protection, EQUOL barely usable). In a sense, he was duplicating much of the effort at Berkeley during the next year, and the two systems quickly and radically diverged.

Issues resolved during this period included:

- 1) updates were defined

- 2) the final syntax and semantics of QUEL were defined
- 3) protection was figured out
- 4) EQUDEL was designed
- 5) concurrency control and recovery loomed on the horizon as big issues. Initial discussions on these subjects started.

2.3 Make It Really Work

The current phase of INGRES development began during the latter part of 1975. At this time the system "more or less" worked. There were lots of bugs and it was increasingly difficult to get them out. The system had performance problems due to convoluted and inefficient code everywhere. The code was also in bad shape. It had been constructed haphazardly by several people, not all of whom were still with the project. Each had his own coding style, way of naming variables, and library of common routines. In short, the system was unmaintainable.

The objective of the current phase was to make the system efficient, reliable, and MAINTAINABLE. At the time we didn't realize that this amounted to a total rewrite. We began to operate with more so-called "controls". There was no more arbitrary tampering with the "current" copy of the code; rudimentary testing procedures were constructed, and rigid coding conventions were enforced. We began to operate

more like a production software house and less like a free wheeling, unstructured operation.

During the current phase concurrency control and recovery were seriously addressed. We took a long time to decide whether to take concurrency control seriously and write a sophisticated locking subsystem (such as the one in SYSTEM-R [GRAY76, GRAY77]) or to do a quick and dirty subsystem using either crude physical locks or predicate locks. We also gave considerable thought to the size of a transaction. Should it be larger than one QUEL statement? If so, the simple strategy of demanding all needed resources in advance and avoiding deadlock was not possible.

The transaction size was eventually decided largely based on simplicity. Once one QUEL statement was selected as the atomic operation for concurrency control and recovery, our hunch was that coarse physical locking would be best. This was later verified by simulation experiments [RIES77, RIES78].

Recovery code was postponed as long as possible because it involved major changes to the utilities. All QUEL statements went through a "deferred update" facility which made recovery from soft crashes (i.e. the disk remains intact) easy if a QUEL statement was being executed. The more difficult problem was to survive crashes while the utilities were running. Each utility performed its own manipulation of

the system catalogs in addition to other functions. Leaving the system catalogs in a consistent state required being able to back up or run forward each command. The basic idea was to create an algorithm which would pass the system catalogs once (or at most twice), find all the inconsistencies regardless of what commands were running, and take appropriate action. Creating such a program required iron clad protocols on how the utilities manipulated the system catalogs. Installing such protocols was a lot of work, most of it in the utilities which everyone by this time regarded as boring code in enormous volume.

The parser had finally become so top heavy from patches that it was rewritten from scratch. Decomposition was improved and the system became progressively faster. In addition, the system was instrumented (no performance hooks were built in from the start). As a result we caught several serious botches. Elaborate tracing facilities were retrofitted to allow a decent debugging environment. In short, the entire system was rewritten.

During this time we also started to support a user community. There are currently some one hundred users -- all requesting better documentation, more features and better performance. These became a serious time drain on the project.

Some of our early users appeared to be contemplating selling

our software. We had taken no initial precautions to safeguard our rights to the code. It became necessary to prepare a licence form and to pull everyone's lawyers into the act. This became a headache that could not easily be deflected, but which made supporting users look easy.

III LESSONS

The following section discusses some of the lessons that have been learned from the INGRES project.

3.1 Goals

Our goals have expanded several times (always when we were in danger of achieving the previous collection). Thus we added features which were not thought about in the initial design (such as concurrency control and recovery) and began worrying about distributed data bases (which was NEVER even talked about earlier). The effect of this goal expansion has been to force us to rewrite a lot of INGRES, in some cases more than once.

3.2 Structured Design

The current wave of structured programming enthusiasts suggests the following implementation plan. Starting with the overall problem, one successively refines it until one has a tree structure of subproblems. Each level in such a tree serves as a "virtual machine" and hides its internal details from higher level machines. We have encountered several

problems in attempting to follow this seemingly sound advice. We discuss four of them.

a) It presumes that one knows what he is doing from the outset. There were many times when we were confused concerning how to proceed. In all cases we chose to do something as opposed to doing nothing, feeling that this was the most appropriate way to discover what we should have done. This philosophy has caused several virtual machines to be dead wrong. Whenever this happened, a lot of redesign was inevitable.

b) We have had to contend with a 64K address space limitation. Initially we did not have a good understanding of how large various modules would be. On more than one occasion we have run out of space in a process which has forced us into the unpleasant task of restructuring the code on space considerations alone. Moreover, since interprocess communication is not fast, we could not always structure code in the "natural" way because of performance problems.

c) There was a strong temptation not to think out all of the details in advance. Because the design leaders had many other responsibilities, we often operated in a mode of "plan the general strategy and rough out the attack". In the subsequent detailed design, flaws would often be uncovered which we had not thought of, and corrective action would have to be taken. Often, major redesigns were the result.

d) It was sometimes necessary to violate the information hiding of the virtual machines for performance reasons. For example, there is a utility which loads indexed sequential (ISAM-like) files and builds the directory structure. It is not reasonable to have the utility create an empty file and then add records one at a time through the access method. This strategy would result in a directory structure with unacceptable performance because of bad balance. Rather, one must sort the records then physically lay them out on the disk and then, as a final step, build the directory. Hence, the program which loads ISAM files must know the physical structure of the ISAM access method. When this structure changed (and it did several times), the loader had to be changed.

All these problems created a virtually constant rewrite/maintenance job of huge magnitude. In four years there have been between two and five incarnations of all pieces of the system. Roughly speaking, we rewrote the majority of the system each year since the project began. Only now is code beginning to have a longer lifetime.

Earlier, there was hesitation on the part of the implementors to document code because it might have a short lifetime. Hence, documentation has been almost non-existent until recently.

3.3 Coding Conventions

To learn the necessity of this task was a very important lesson to us. As mentioned earlier, the equivalent of one total rewrite resulted from our initial failure in this area. We found that pieces of code which had a non trivial lifetime were unmaintainable except by the original writer. Also, every time we gave someone responsibility for a new module he or she would rewrite it according to his or her standards (allegedly to clean up the other person's bad habits). This process never converges and I feel that it is similar to the dog or wolf who stakes out his "turf" by urinating on each bush on its perimeter.

Only coding conventions stop this process.

3.4 User Support

There are lessons which we have learned about users in three areas.

3.4.1 Serious Users

There are a few serious users (5-10). All have been extremely bold and forward-looking people and have exercised our system extensively before committing to use it. All of these users first chose UNIX (which says something about their not being a random sample of users) and then obtained INGRES.

Most have made modifications to personalize INGRES to their needs, viewed us as a collection of goofy academicians and

were pretty skeptical that our code was any good. All were very concerned about support, future enhancements and how much longer our research grants would last.

All have developed end user facilities using EQUEL and have given us a substantial wish list of features. The following is typical:

- 1) the system is too slow (especially for trivial interactions)
- 2) the system is too slow for very large data bases (whatever this means)
- 3) protection, integrity constraints and concurrency control are missing (true for earlier versions)
- 4) the EQUEL interface is not particularly friendly
- 5) the system should have partial string matching capabilities, a data type of "bit", and a macro facility. (The wish list of such features is almost unbounded.)

Surprisingly, nobody has ever complained about the crash recovery facilities. Also, a concurrency control scheme consisting of locking the whole data base would be an acceptable alternative for most of our users.

The biggest problem that these users have faced is the problem of understanding some 400,000 bytes of source code, most of it free of documentation (other than comments in the

code).

The merits of INGRES that most of these users claim, rest on

1) ease of use. The system is easy to use after a minor amount of training. The "startup" cost is much lower than for other systems.

2) The high level language allows applications to be constructed incredibly fast, as much as 10 times faster than originally anticipated.

This short coding cycle allowed at least one user to utilize a novel approach to application design. The conventional approach is to construct a specification of the application by interacting with the end user. Then programmers go into their corner to implement the specifications. A long time later they emerge with a system and the users respond that it is not really what they wanted. Then, the rounds of retrofitting begin.

The novel approach was to do application specification and coding in parallel. In other words, the application designer interacted with end users to ascertain their needs and then coded what they wanted. In a few days he returned with a working prototype (which of course was not quite what they had in mind). Then the design cycle iterated. The important point is that end users were in the design loop and their needs were met in the design process. Only the

ability to write data base applications quickly and economically allowed this to happen.

3.4.2 Casual Users

There are about 90 more "casual" users. We hear less from these people. Most are universities who use the system in teaching and research applications. These users are less disgruntled with performance and unconcerned about support.

3.4.3 Performance Decisions

Users are not always able to make crucial performance decisions correctly. For example, the INGRES system catalogs are accessed very frequently and in a predictable way. There are clear instructions concerning how the system catalogs should be physically structured. Even so, some users fail to make the necessary modifications. Of course, the system continues to run, it just gets slower and slower. Finally, we removed this particular decision from the users domain entirely. It makes me a believer in automatic data base design (e.g. [HAMM76])!

IV FLAT OUT MISTAKES

This section will discuss what we believe to be the major mistakes in the current implementation.

4.1 Interpreted Code

The current prototype interprets QUEL statements even when

these statements come from a host language program. An interpreter is reasonable when executing ad-hoc interactions. However, the EQUQL interface processes interactions from a host language program as if they were ad-hoc statements. Hence, parsing and finding an execution strategy are done at run time, interaction by interaction.

The problem is that most interactions from host languages are simple and are done repetitively. (For example, giving a 10 percent raise to a collection of employee names read in from a terminal amounts to a single parameterized update inside a WHILE statement). The current prototype has a fixed overhead per interaction of about 400 msec. Hence, throughput for simple statements is limited by this fixed overhead to about 2.5 interactions per second. Parsing at compile time would reduce this fixed overhead somewhat.

At least as serious is the fact that the interpreter consumes a lot of space. The "working set" for an EQUQL program is about 150K bytes plus the program. For systems with a limited amount of main memory this presents a terrible burden. A compiled EQUQL would take up much less space (at least for EQUQL programs with fewer than 10 interactions per program). Moreover, a compiled EQUQL could run as less processes, saving us some interprocess communication overhead. This issue is further discussed in Section 4.3.

The interpreter was built with the notion of ad-hoc interac-

tions in mind. Only recently did we realize the importance of a programming language interface. Now we are slowly converting INGRES to be alternatively compiled and interpreted. We were clearly naive in this respect.

4.2 Validity Checking

This mistake is related to the previous one. When an interaction is received from a terminal or an application program, it is parsed at run time. Moreover, (and at a very high cost) the system catalogs are interrogated to validate that the relation exists, that the domains exist, that the constants to which the domains are being compared are of the correct type or are converted correctly, etc. This costs perhaps 100 msec. of the 400 msec. fixed overhead, and no effort has been made to minimize its impact. This makes the "do nothing" overhead high and, from a performance viewpoint, is the really expensive component of interpretation.

4.3 Process problems

The "do nothing" overhead is greatly enlarged by our problems with a 16 bit address space. The current system runs as 5 processes (and the experimental system at Berkeley as 6) and processing the "nothing" interaction requires that the flow of control go through 8 processes. This necessitates formatting 8 messages, calling the UNIX scheduler 8 times and invoking the interprocess message system (pipes) 8

times. This generates about 150-175 msec. of the 400 msec. of fixed overhead.

In addition, code cannot be shared between processes. Hence, the access methods must appear in every process. This causes wasted space and duplicated code. Moreover, a UNIX file can only be opened by one process on behalf of itself. Since each process must look at the system catalogs they must be opened individually by each process. Again there is considerable repetition.

Besides this performance problem, the previous section noted that the process structure has changed several times because of space considerations. As a result, a considerable amount of energy has gone into designing new process structures, writing the code which correctly "spawns" the right run time environment and handling user interrupts correctly.

In retrospect, we had no idea how serious the performance problems associated with being forced to run multiple processes would be. It would have been clearly advantageous to choose a 32 bit machine for development; however, there was no affordable candidate to be obtained at the time we started. Also, perhaps we should have relaxed the 64K address limitation once we obtained a PDP-11/70 (which has a 128K limitation). This would have cut the number of processes somewhat. However, many of our 100 users have 11/34's or 11/40's and we were reluctant to cut them off.

Lastly, we could have opted for less complexity in the code. However, to be effective, the system would have to be reduced by at least a factor of two. It is not clear that an interesting system could be written within such a constraint. The bottom line is that this has been an enormous problem, but one for which we see no obvious solution, other than to buy a PDP-11/780 and correct the situation now that a 32 bit machine exists which can run our existing code.

4.4 Access Methods

The decision was made very early that we were not going to write our own file system to get around UNIX performance (as SYSTEM-R elected to do [ASTR76]). Instead, we would simply build access methods on top of the existing file system.

The reasoning behind this decision was to avoid duplicating operating system functions. Also, exporting our code would have been more difficult if it contained its own file system. Lastly, we underestimated the severity of the performance degradation that the UNIX file system contributes to INGRES when it is processing large queries. This topic is further discussed in [HAWT79]. In retrospect, we probably should have written our own file system.

The other problem with the access methods concerns whether they are I/O bound. Our initial assumption was that it would never take INGRES more than 30 msec. to process a 512 byte page. Since it takes UNIX about this long to fetch a

page from the disk, INGRES would always be I/O bound for systems with a single disk controller (the usual case for PDP-11 environments). Although INGRES is sometimes I/O bound, there are significant cases where it is CPU bound [HAWT79].

The following three situations are bad mistakes when INGRES is CPU bound:

- a) An entire 512 page is always searched even if one is looking only for one tuple (e.g. a hash bucket is a UNIX page).
- b) A tuple may be moved in core one more time than is strictly necessary.
- c) A whole tuple is manipulated rather than just desired fields.

Although we have corrected points b) and c), point a) is fundamental to our design and is a mistake.

4.5 Static Directories

INGRES currently supports an indexing access method with a directory structure which is built at load time and never modified thereafter. The arguments in favor of such a structure are presented in [HELD78]. However, we would implement a dynamic directory (as in B-trees) if the decision were made again. Two considerations have influenced

the change in our thinking.

The data base administrator has the added burden of periodically rebuilding a static directory structure. Also, he can achieve better performance if he indicates to INGRES a good choice for how full to load data pages initially. In the previous section we indicated that data base administrators often had trouble with performance decisions, and we now believe that they should be relieved of all possible choices. Dynamic directories do not require periodic maintenance.

The second fundamental problem with static directories is that buffer requirements are not predictable. In order to achieve good performance, INGRES buffers file system pages in user space when advantageous. However, when overflow pages are present in a static directory structure, INGRES should buffer all of them. Since, address space is so limited, a fixed buffer size is used and performance degrades severely when it is not large enough to hold all overflow pages. On the other hand, dynamic directories have known (and nearly constant) buffering requirements.

4.6 Decomposition

Although decomposition [WONG76] is an elegant way to process queries which is easy to implement and optimize, there is one important case which it cannot handle. For a two variable query involving an equi-join, it is sometimes best to

sort both relations on the join field and then merge the results to identify qualifying tuples [BLAS77]. It is impossible for us to add this as a tactic and apply it when it is appropriate without dramatically altering the INGRES process structure. Again, the address space issue rears its ugly head!

4.7 Protection

It appears much cleaner to protect "views" as in [GRIF76] rather than base relations as in [STON74, STON76]. It appears that sheer dogma on my part prevented us from correcting this.

4.8 Lawyers

I would be strongly tempted to put INGRES into the public domain and delete our interactions with all attorneys (ours and everyone else's). Whatever revenue the University of California derives from license fees may well not compensate for the extreme hassle which licencing has caused us. Great insecurity and our egos drove us to force others to recognize our legal position. This was probably a big mistake.

4.9 Useability

Insufficient attention has been paid to the INGRES user interface. We have learned much about "human factors" during the project and have corrected many of the botches. However, there are several which remain. Perhaps the most

inconvenient is that updates are "silent". In other words, INGRES performs an update and then responds a "done". It never gives an indication of the tuples that were modified, added or deleted (or even how many there were). This "feature" has been soundly criticized by almost everyone.

V COMMENTS

This section contains a collection of comments about various things which do not fit easily into the earlier sections.

5.1 UNIX

As a program development tool, we feel that UNIX has few equals. We especially like the notion of the command processor, the notion of pipes, the ability to treat pipes, terminals and files interchangeably, the ability to spawn subprocesses and the ability to fork the command interpreter as a subprocess from within a user program. UNIX supports these features with a pleasing syntax, very few "surprises", and most unnecessary details (e.g. blocking factors for the file system) remain hidden.

The use of UNIX has certainly expedited our project immeasurably. Hence, we would certainly choose it again as a operating system.

The problems which we have encountered with UNIX have almost all been associated with the fact that it was envisioned as a general purpose time-sharing system for small machines and

not as a support system for data base applications.

Hence, there is no concurrency control and no crash recovery for the file system. Also, the file system does not support large files (16 Mbytes is the current limit) and uses a small (512 bytes) page size. Moreover, the method used to map logical pages to physical ones is not very efficient. In general, it appears that the performance of the file system for our application could be dramatically improved.

5.2 The PDP-11

Other than the address space problems with a PDP-11, I have only two other comments regarding the hardware. First, there is no notion of "undefined" as a value for numeric data types supported by the hardware. Allowing such a notion in INGRES would require taking some legal bit pattern and by fiat making it equal undefined. Then we would have to inspect every arithmetic operation to see if the chosen pattern happened inadvertantly. This could be avoided by simple hardware support (such as found on CDC 6000 machines).

Second there is no machine instruction which can move a string in core. Consequently, data pages are moved in core one word at a time inside a loop. This is a source of considerable inefficiency.

5.3 Data Models

There has been a lot of debate over the efficiency of the various data models. In fact, a major criticism of the relational model has been its (alleged) inefficiency.

There are (at least) two ways to compare the performance of data base systems.

a) the overhead for small transactions. This is a reasonable measure for how many transactions per second can be done in a typical commercial environment.

b) The cost of a given big query

It should be evident that a) has nothing to do with the data model used (at least in a PDP-11 environment). It is totally an issue of the cost of the operating system, system calls, environment switches, data validity costs, etc. In fact, if INGRES were a network oriented system and ran as five processes, it would also execute 2.5 transactions per second.

The cost of a big query is somewhat data model dependent. However, even here this cost is extremely sensitive to the cost of a system call, the operating system decisions concerning buffering and scheduling, the cost of shuffling output around and formatting it for printing, and the extent to which clever tuning has been done. In addition, the design of a data base management system is often very sensitive to the features (and quirks) of the operating system on which

it is constructed. (At least INGRES is). These are probably much more important in determining performance than what data model is used.

In summary, I would allege that a comparison of two systems using different data models would result primarily in a test of the underlying operating system and the implementation skill (or man years allowed) of the designers and only secondarily in a test of the data models.

VI INGRES PROJECT PLANS

INGRES appears to be at least potentially commercially viable. However, a commercial version would require, at least:

- 1) someone to market it
- 2) much better documentation
- 3) someone willing to guarantee maintenance. (Whether or not we do it, the University of California will not promise to fix bugs.)
- 4) a pile of boring utilities (e.g. a report generator, a tie in to some communications facilities, and access to the system from other languages than C).

Even so, we would not have a good competitive position because UNIX is not supported and because no COBOL exists for UNIX.

There has been a clear decision on the part of the major participants not to create a commercial product. On the other hand, the project cannot simply announce that it has accomplished its goals and close shop. Hence, we have gone through a (sometimes painful) process of self examination to decide "what next". Here are our current plans.

1) distributed INGRES

We are well into designing a distributed data base version of INGRES which will run on a network of PDP-11's. The idea here is to hide the details of location of data from the users and fool them into thinking that a large unified data base system exists [EPST78, STON78].

2) A distributed data base machine

This is a variant on a distributed data base system in which we attempt only to improve performance. It has points in common with "back end machines" and amounts to moving code from a host UNIX into multiple slave "back ends" [STON78a].

3) A new data base programming language

Obviously starting with C and an existing data base language QUEL and attempting to glue them together into a composite language is rather like interfacing an apple to a pancake. It would clearly be desirable to start from scratch and design a good language. Initial thoughts on this language are presented in [PREN77].

4) A data entry facility

Now that the component of writing transactions which can be attributed to the data base system has shrunk to near zero (by high level language facilities), we are left with transactions that have virtually no data base code and are entirely what might be called "screen definition, formatting and data entry". We are designing a facility to help in this area.

5) Improved integrity control

Currently, INGRES is not very smart in this area. Other than integrity constraints [STON75] (which do something but not as much as might be desired), we have no systematic means to assist users with integrity/validation problems. We are investigating what can be done in this area.

It is pretty clear that all of the above will require substantial changes in the current software. Hence, we can remain busy for a seemingly arbitrary amount of time. This will clearly continue until we get tired or are again in danger of meeting our goals.

ACKNOWLEDGEMENT

The INGRES project has been directed by Professors Eugene Wong and Larry Rowe in addition to myself. The role of chief programmer has been filled by Gerald Held, Peter

Kreps, Eric Allman and Robert Epstein. The following persons worked on the project at various times; Richard Berman, Ken Birman, James Ford, Paula Hawthorn, Nancy MacDonald, Daniel Ries, Peter Rubinstein, Michael Ubell, Nick Whyte, Carol Williams, Karel Yousseffi and William Zook.

The INGRES project is sponsored by the U.S. Air Force Office of Scientific Research Grant 78-3596, the U.S. Army Research Office Grant DAAG29-76-G-0245, the Naval Electronics Systems Command Contract N00039-78-G-0013 and the National Science Foundation Grant MCS75-03839-A01.

REFERENCES

- [ALLM76] Allman, E., Held, G. and Stonebraker, M., "Embedding a Data Manipulation Language in a General Purpose Programming Language," Proc. 1976 ACM-SIGPLAN-SIGMOD Conference on Data Abstractions, Salt Lake City, Utah, March, 1976.
- [ASTR76] Astrahan, M. M. et. al., "System R: A Relational Approach to Database Management," TODS 2, 2, June 1976.
- [BLAS77] Blasgen, M. and Eswaren, K., "Storage and Access in Relational Data Base Systems," IBM Systems Journal, December, 1977.

- [EPST78] Epstein, R., Stonebraker, M., and Wong, E., "Query Processing in a Distributed Data Base System," Proc. 1978 ACM-SIGMOD Conference on Management of Data, Austin, Texas, May, 1978.
- [GRAY76] Gray J. et. al., "Granularity of Locks and Degrees of Consistency in a Shared Data Base," IBM Research, San Jose, Ca., RJ 1849, July, 1976.
- [GRAY77] Gray, J., "Notes on Data Base Operating Systems," unpublished course notes, July 1977.
- [GRIF76] Griffiths, P. and Wade, B., "An Authorization Mechanism for a Relational Data Base System," TODS, 2, 3, September 1976.
- [HAMM76] Hammer. M. and Chan, I., "Index Selection in a Self Adaptive Data Base System," Proc. 1976 ACM-SIGMOD Annual Conference on Management of Data, Washington, D.C., June 1976.
- [HELD78] Held, G. and Stonebraker, M., "B-Trees Reexamined," CACM, February, 1978.
- [HAWT78] Hawthorn, P. and Stonebraker, M., "Use of Technological Advances to Enhance Data Base Management System Performance," Electronics Research Laboratory, University of California, Memo No. 79-5, January, 1979.
- [PREN77] Prenner, C. and Rowe, L., "Programming Languages for Relational Data Base Systems," Proc. 1978

National Computer Conference, Anaheim, Ca., June, 1978.

- [RIES77] Ries, D. and Stonebraker, M., "A Study of the Effect of Locking Granularity in a Relational Data Base System," TODS 3, 3, September 1977.
- [RIES78] Ries, D. and Stonebraker, M., "Lock Granularity Revisited," to appear in TODS.
- [RITC75] Ritchie, D. and Thompson, K., "The UNIX Time-Sharing System," CACM, June 1975.
- [STON74] Stonebraker, M. and Wong, E., "Access Control in a Relational Data Base System by Query Modification," Proc. 1974 ACM Annual Conference, San Diego, Ca., November 1974.
- [STON75] Stonebraker, M., "Implementation of Integrity Constraints and Views by Query Modification," Proc. 1975 ACM-SIGMOD Conference on Management of Data, San Jose, Ca., June 1975.
- [STON76] Stonebraker, M. and Rubinstein, P., "The INGRES Protection System," Proc. 1976 ACM Annual Conference, Houston, Texas, November 1976.
- [STON76a] Stonebraker, M. et. al., "The Design and Implementation of INGRES," TODS 2, 3, September 1976.
- [STON78] Stonebraker, M., "Concurrency Control, Crash Recovery and Consistency of Multiple Copies of Data in a Distributed Data Base System," Proc. 3rd

Berkeley Workshop on Distributed Data Bases and
Computer Networks, San Francisco, Ca., August,
1978.

[STON78a] Stonebraker, M., "A Distributed Data Base
Machine," Electronic Research Laboratory, Univer-
sity of California, Memo No. M78-55, June 1978

[WONG76] Wong, E. and Youseffi, K., "Decomposition: A Stra-
tegy for Query Processing," TODS, 2, 3, September
1976.