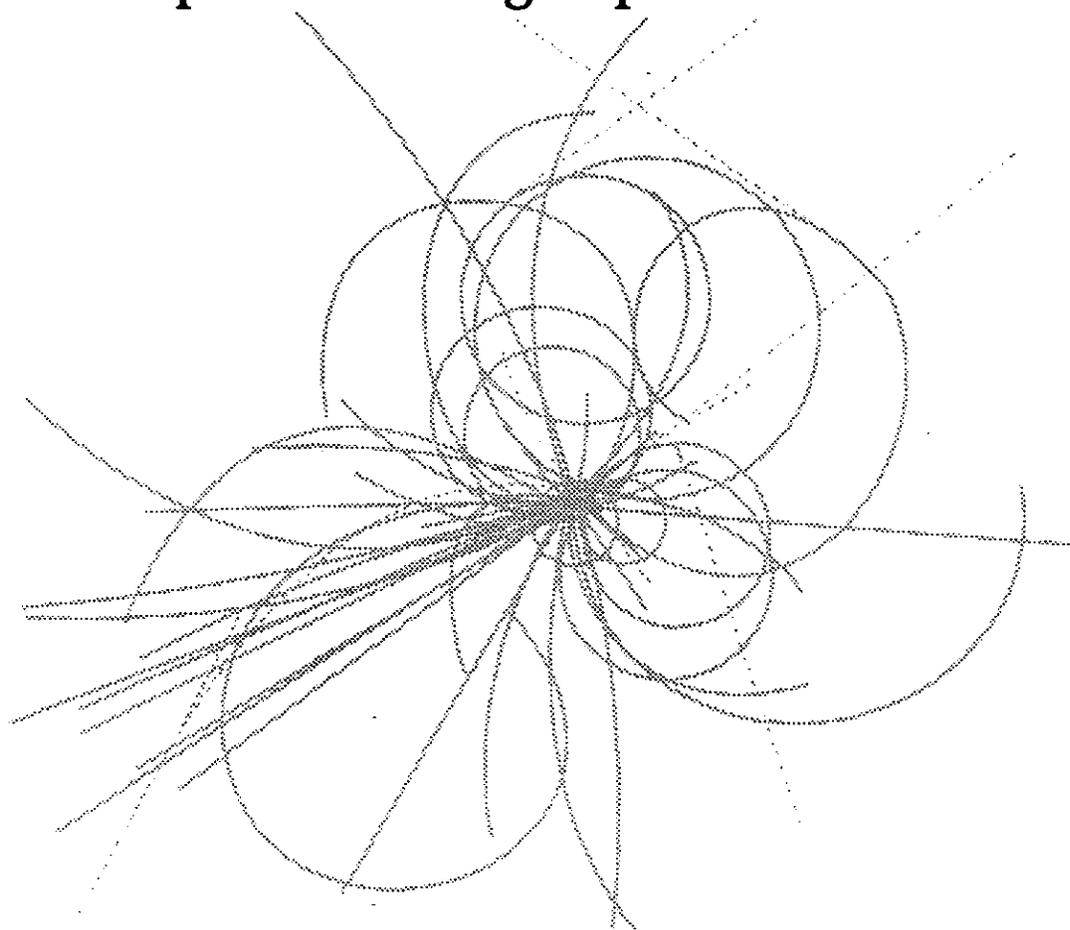


Superconducting Super Collider Laboratory



A Model for Computing at the SSC

Drew Baden and Robert Grossman

June 1990

A MODEL FOR COMPUTING AT THE SSC

Drew Baden

Department of Physics, University of Maryland

Robert Grossman

Laboratory for Advanced Computing

Department of Mathematics, Statistics, and Computer Science

University of Illinois, Chicago

June 1990

TABLE OF CONTENTS

	Title	Page
1.	Introduction	
2.	HEP: Computing	
2.1	ONLINE	
2.1.1	Data Storage Media	
2.1.2	Online Storage Proposal	
2.2	PRODUCTION (of DSTs)	
2.2.1	Proposal	
2.3	ANALYSIS	
2.3.1	Data Structure	
2.3.2	The Relational Database	
2.3.3	The Extensible Database	
2.3.4	The Scientific Database	
2.3.5	A HEP Database	
2.3.6	Building a HEP Database	
3.	Remote Access: Client/Server	
3.1	Collaborating University/Laboratory Access	
4.	SSCL Implications	
4.1	SSCL Computer System	
4.2	Existing Resources	
4.3	Key Factors	

LIST OF TABLES

1.	KODAK vs. SUMMUS r/w optical disk comparison	
2.	Data storage media comparisons	

A MODEL FOR COMPUTING AT THE SSC

Drew Baden

Department of Physics, University of Maryland

Robert Grossman

Laboratory for Advanced Computing

Department of Mathematics, Statistics, and Computer Science

University of Illinois, Chicago

June 1990

ABSTRACT

High energy physics experiments at the Superconducting Super Collider (SSC)* will show a substantial increase in complexity and cost over existing forefront experiments, and computing needs may no longer be met via simple extrapolations from the previous experiments. We propose a model for computing at the SSC based on technologies common in private industry involving both hardware and software.

* Operated by Universities Research Association, Inc., for the U.S. Department of Energy under Contract No. DE-AC02-89ER40486.

1. INTRODUCTION

High energy physicists have traditionally had a schizophrenic attitude towards computer hardware and software. Demand for increases in computer performance (CPU, I/O, data acquisition, etc.) in high energy physics (HEP) has consistently increased over time, and contributes to some degree, to the evolution of computers. On the other hand, the language of choice has always been Fortran (considered extremely cumbersome for some applications). Physicists in general are reluctant to take time to learn how the computer really works. Software is considered an unworthy thing to work on. There is justification for some attitudes. In order to be effective, the physicist must be very comfortable with the computer and be able to devote all of their attention to the problem the physicist is trying to solve. However, the increasing demands the physicist makes (and will continue to make at the SSC) in computing, may necessitate a basic change in the way things are done. A novel opportunity is available to review computing at the SSC for two reasons:

1. There are no archaic structures already set up (which is a basic cause of computer-conservatism, or inertia).
2. The amount of data both online and offline anticipated at the SSC may already be too much to handle with present "classical" HEP computing techniques.

This paper outlines a scheme for computing at the SSC with an emphasis on data analysis. Such computing must incorporate state-of-the-art computer science tailored to the needs of HEP.

2. HEP: COMPUTING

Data analysis in HEP has traditionally consisted of the following:

1. **ONLINE:** Collect data, write to tape.
2. **PRODUCTION:** Create "Data Summary Tapes" (DSTs) by reading from tape and running "production" code(s) which reconstruct the data into dimensional quantities. Results are written to some medium, usually magnetic tape.
3. **ANALYSIS** (of DSTs): This stage consists of performing calculations and making cuts based on either calculation results or variables already present in the record. The entire record (event) is read in for each event analysis.

Historically, each time an experiment increases in complication, the weak parts of the above outline have traditionally been dealt with through a logical extension of traditional techniques, incorporating more technology. It is possible this can continue, given the impressive technological advances in present and anticipated computer hardware. However, the future (SSC) will most likely result in an order of magnitude increase in the size of each event (relative to the largest ongoing HEP experiments to date), the number of events needed for analysis, and in the number of physicists who will be placing demands on a particular system (the size of the SDC is expected to approach $\simeq 1000$ by the time data is collected). This amounts to a large, non-linear increase in the computing requirements. The question is, can we continue to throw all available increases in technology using the traditional techniques at our computing problems? If

not, can an advantage be taken of both long large lead times now available before the data is collected and the tremendous increases in data storage, CPU performance, and software designs either present or anticipated to redefine the way computing is accomplished, such that technological developments will help solve problems? The safest thing to do is to anticipate that we will have to do a bit of rethinking about the way we do computing. This means that we will first have to educate ourselves about what is currently available in the way of computer science (software and organization/management of resources) which might be applicable to HEP needs. In order to do this, a rationale scheme must be developed, taking into account the ramifications of the three categories, **ONLINE**, **PRODUCTION**, and **ANALYSIS** listed above.

2.1 ONLINE

During the 1988–1989 Tevatron run, the CDF experiment wrote 80–100 kbyte event records to 9 track tape at approximately 1 Hz. Each tape can hold approximately 120 Mbytes at 6250 bpi, and with a 20% safety factor this limited the tapes to about 1000 events each. At 1 Hz, a tape was filled every 15 minutes.

Near the end of the run, the data was stored on 8 mm cassettes, which can hold approximately 2.3 Gbytes, a factor of about 20 greater than 9 track tapes. The 8 mm drives can write data at about 250 kbyte/sec. This seems satisfactory for CDF needs, providing that the product of the event size times the event rate never exceeds $\# \text{tapes} \times 250 \text{ kbyte/sec}$ given today's technology. EXABYTE (the 8 mm of choice at FNAL) has just started to ship the next generation of 8 mm drives. Improvements include:

- Factor of 2 increase in bandwidth to 500 kbyte/sec
- Factor of 2 increase in capacity to 4.6 Gbytes

However, 8 mm tapes can only be seen as a solution to long-term storage problems. Access to data is quite limited on an interactive basis, since the device is not random-access. A simple calculation based on the above numbers say that it takes $\simeq 2.5$ hours to read an entire tape. This is not as much of a problem for the **ONLINE** storage as it is for the **PRODUCTION** and **ANALYSIS**, however it does not make sense to support too many different storage media at SSCL. In Texas, 8 mm tapes will probably not be sufficient. Anticipated event sizes of the order of 1 Mbyte would require 20 8 mm tapes online to get a rate of 10 Hz (given the writing speed is 500 kbyte/sec).

NOTE

These tapes are guaranteed to write with an error rate of less than 1 in 10^{15} bits and read with an error rate of less than 1 in 10^{13} bits. Given a 1 Mbyte event size, this corresponds to a write error of less than 1 in 125 million events and a read error of less than 1 in 1.25 million events.

Data Storage Media: Optical Disks

At the SSC, optical disks will have to be considered. Such devices are relatively new on the market, but have already reached milestones in performance worth considering. Write-Once-Read-Many (WORM) drives have been available for several years from various manufacturers.

In late 1988, erasable optical drives started to become available. Although slower than conventional magnetic disks, erasable optical drives have the advan-

tage of a larger information density (currently a factor of 2-10 over magnetic) and removable media.

For example, KODAK makes a device called the "OPTICAL DISK SYSTEM 6800" and SUMMUS now markets a similar device called the "LIGHTDISK SO-600". The KODAK device has a larger capacity but is slower and more expensive than the SUMMUS device. See Table 2.1.1 for a comparison.

Table 2.1.1 Comparison Between KODAK's 6800
and SUMMUS's LIGHTDISK

Characteristics	KODAK	SUMMUS
Bus	SCSI,IPI-3	Unibus, Q-BUS,VME,SCSI
Capacity	6.8 Gbyte (13.2MBlocks)	594 Mbyte (1.2MBlocks)
Volume (in ³)	3000	150
Weight	27 kg	3 kg
Transfer Rate	1 Mbyte/sec	680 kbyte/sec
Access Time	700 ms (Data within disk)	90 ms (average seek)
Price DRIVE	\$20,000	\$6,000
Price DISK	\$700	\$250

The KODAK device can be purchased with a robotic library server (called a "jukebox") which holds 250 disks with an advertised access time of 6.5 seconds per disk costing from \$100,000 to \$200,000. SUMMUS has not announced such a product, but it will probably be available some time in the near future.

Data Storage Media: Optical Tape

CREO PRODUCTS (Vancouver, British Columbia) has announced the availability some time this year of a new optical tape drive, using the so-called "optical paper" development by ICI. Honeywell Test Instruments Division has reportedly agreed to market the drive to end users. The characteristics of this drive are

- 12 inch diameter tapes, 1/2 inch wide (similar to conventional 9-track tapes) and 2400 feet long
- 1 TByte capacity per tape
- 3 Mbyte/sec transfer rate
- 28 sec average access time to any byte on any tape, with 60 sec. end-to-end
- \$225,000 per drive
- \$5000-10,000 per tape
- Write-once media.

One of the biggest advantages of having such a medium comes from its extremely large capacity. A terabyte is approximately 10^6 events. This entails 100 tapes to store an eventual $\simeq 10^8$ events, which means a large cost in tapes given the above price. Should the price drop a factor of 10, it would no longer be a consideration. However, this solution is basically just a technological extension of 8 mm tapes in that it is really not a device that one could use for interactive computing. This issue is discussed further below.

Data Storage Media: RAID

Called Redundant Array of Inexpensive Disks, RAID is a scheme which uses a farm of small disks with multiple parallel data paths in order to simulate the

far higher performance of a large disk farm. In the limit of an infinite number of disks and data paths, the fastest way to write data is to write 1 bit-per-disk. The RAID project uses block transfers of order $\simeq 4$ kbits writing over parallel links to 500 Mbyte disks, and is used by multi-processor farms (connection machine-type computers). It is worth investigating, but the small capacity compared to optical disks means that a great number of disk drives will have to be purchased and managed. However, assume that in $\simeq 5$ years time these disks will have a 1 Gbyte capacity, 1 Mbyte/sec bandwidth, and \$1000 price tag. If events are written at 100 Hz and are 1 Mbyte long, then at least 200 disks would be required for online purposes. To store 10^7 events of 1 Mbyte length, it would require 10 000 disks at a cost of \simeq \$10M. The cost would probably eliminate this as a choice for mass storage, but may play a role in caching.

Online Storage Proposal

Table 2.1.2 provides a comparison between the SUMMUS 8 mm Gigatape (today's standard) using the new EXABYTE specifications (0.5 Mbyte/sec, 4.6 Gbyte storage), the KODAK read/write optical disk (as an alternative device which has random access capabilities) and the CREO optical tape (for data storage). Events are assumed to be 1 Mbyte long, and it is assumed that $\simeq 10^7$ events will be collected during the first year's running.

Table 2.1.2 Comparison between various optical disk, optical tape, and 8 mm magnetic tape, and magnetic disks for online and offline storage. Event size is assumed to be 1 Mbyte and number of events stored is estimated to be 10^7 for the first year of running

Device	Capacity (1 MB events)	Write Speed (events/sec)	Cost (10^7 events)	Interactive?
Optical Disk (KODAK)	6800	1	\simeq \$1.2M	YES
8 mm tape (SUMMUS)	4600	0.5	\simeq \$0.2M	NO
Optical Tape (CREO)	1 000 000	3	\simeq \$1.0M	NO
RAID	10 000 000	≥ 100	\simeq \$10.0M	YES

Given the online specifications of an SSC experiment (assuming 1 Mbyte events, 10 Hz to tape), a proposal for online storage using the KODAK device would entail:

- Drives: At 10 Hz, ten of the KODAK-type optical disk drives. The cost of this before university discount will be on the order of \simeq \$200,000. Each drive's disk would have to be changed every $\simeq 2$ hours, requiring a disk change somewhere every $\simeq 12$ minutes.
- Disks: Each disk holds 6800 events. A total of $\simeq 10^7$ events would require $\simeq 1500$ disks for the online. The cost of this (at \simeq \$700/disk) is \simeq \$1,000,000 before university discount.

It can be assumed that these prices are upper limits in relation to the eventual real cost at the time the experiments are running.

If the data taking were to begin today, these costs would in no way be prohibitive, relative to a cost of \simeq \$500M for each SSC experiment. However, the real power using random access devices for mass storage is realized in the advantages for the other two categories of HEP computing, **PRODUCTION** and **ANALYSIS**, is outlined below. It is not *necessary* that data be written online, onto a random access device, except it will be desirable to limit the different number of storage-for-support reasons to as few as possible, and given that one will want to minimize the number of times one is copying such large amounts of data from one device to another.

2.2 PRODUCTION (OF DSTs)

There are two scenarios to consider for how offline production will proceed:

1. High-level objects (electrons, muons, jets, photons, etc.) are “produced” offline using special-purpose codes operating directly on the raw data. The program will use calibration constants which are calculated using the raw data itself after (or as) it is written.
2. Production will be accomplished in the last trigger level (Level 2) of the DAQ system. This is possible only if the final constants are available in real time and there is enough computing power available.

The choice will depend on the type of hardware which will be used in the detector. The safest choice would be to prepare for item 1.

Proposal

Assuming item 1 above, one could imagine that a farm of processors, with access to the database of constants and the data via an optical “juke-box” (see 2.1.1 Data Storage Media) would suffice. Note that unlike 8 mm tapes, there is no need for staging. The numbers look something like this:

1. On the order of 2 “juke-boxes”, one for providing data to be read by the production program and one for storing the resultant DSTs. Each “juke-box” is priced at \simeq \$100,000 to \$200,000. Total cost is no more than \simeq \$400,000 at today’s prices before university discount.
2. The second “juke-box” will require additional optical disks (relative to the ones used to record the data online). One can probably come fairly close to reality by estimating that the number of disks used by DSTs is no larger than approximately two times that used for the raw data (assuming an expansion of 2 and no reduction in the number of events). The cost of this, again at today’s prices before discount, is \simeq \$1,000,000 for 10^7 events. To speed up access to the DSTs, the juke-box of optical disks should be provided with magnetic disk caching. Hybrid systems of this type are now becoming commercially available. Such systems could be used off-the-shelf or higher performance and more specialized caching systems could be designed and developed expressly for this purpose.

2.3 ANALYSIS

Traditionally, each new HEP experiment has invented its own software “system”. Physicists then use this system for **PRODUCTION** and **ANALYSIS**. This has always been simply a matter of convenience. At the SSC, such “convenience” may never appear, given $\simeq 1000$ physicists and $\simeq 10^7$ events with $\simeq 1$ Mbyte per event, and one has to consider a different “system” for **PRODUCTION** and **ANALYSIS**. Experience at CDF reinforces this for a number of reasons:

1. *Linking*. Most CDF analysis is VAX-based. The offline software packages are quite large, and even with shareable (and shadowing) link times are not in the “interactive” ballpark range. Even if they were, the nature of the physicist ensures that (almost) no code freshly written will ever run correctly the first time, which means that a non-trivial fraction of the computer resources are devoted to linking.
2. *Data*. Data at CDF is written in a device-independent manner, necessitated by the device dependencies of the collaboration. This means one has to add structure to the data, and this increases its length and complexity. Given that CDF is dependent on magnetic tape for mass storage and magnetic disks for interactive storage, data space for the individual user is at a premium.

Most of the CDF physicists are turning to the CERN-written PAW program to produce physics results. This program allows interactive data analysis (much like a spreadsheet program) and incorporates an inline Fortran interpreter which

completely eliminates the need for linking. The interpreter has been found to be fast enough, since most of the time spent when running in the interpreter mode is for input/output (I/O). The scheme for data analysis looks like this:

1. Binary files (device-independent, unformatted I/O via very simple Fortran `WRITE(LUN) x,y,z...`) are made from the DSTs. These files contain only specific parts of each event necessary to perform the particular analysis as determined by each physicist. The files are created by running a single job using the **PRODUCTION** “system”, and the data are stored on disk as reals and integers as appropriate. These files are extremely small relative to the DSTs (factor of at least 100) and easy to keep track of.
2. **HBOOK 4 (PAW)** files are created by running a very simple interactive program which reads the binary data and fills “ntuples” (a matrix where the rows are events and the columns are physics quantities). There is very little to this program other than straight Fortran and some **HBOOK** commands with the proviso that care must be taken to keep track of the words in the binary file, since it is not self-describing. The PAW file consists of **ZEBRA** banks, which is not in any way compact, but can be remade quickly and easily, provided the binary file is available. This (PAW) file is needed to run the PAW program.
3. Run PAW. The program provides the ability to make histograms, scatter plots, fits (interactive **MINUIT**), and allows one to place cuts on the data and study correlations and anti-correlations interactively. PAW is not easy to use—the commands are “UNIX”-based and the documentation has

been found to be less than clear—but the program is extremely flexible and powerful and has become quite popular at CDF.

The popularity and success (given that CDF physics output is not unsatisfactory) of the PAW program warrants consideration at the SSC. However, the dependency on binary files is bothersome for a number of reasons (e.g. files are device dependent and unstructured hence not easily self-describing), and the file structure of the PAW file is very primitive. It is conceptually straightforward to update the concept.

Data Structure

One can easily improve on the above “PAW” scenario to minimize waste of resources. First, consider the physicist who is either resident at SSCL or can log into a central computer system over a high speed link (the kind of T1 lines which are proliferating now). Given the above proposal to write the data to optical disks [see Proposal (2.2.1)] large amounts of data are accessible at the SSCL without having to go to tape (the “juke-box” allows access to about 1 Tbyte of data, or about 10^6 events on disk). Given the belief that it will be I/O and not CPU which will limit us, we can consider ways of restructuring the way data is stored and analyzed so as to minimize the I/O. If one accepts that there may be a benefit in separating computing in HEP into **PRODUCTION** and **ANALYSIS** problems, then any and all solutions to (software) computing issues will have to start with the way the data is structured. Data structures should be designed to maximize the benefits of the increasing technology—the two are not independent. Intelligent data structures could greatly aid in bookkeeping efforts

during production as well as reduce the terrific amount of I/O required for data analysis.

Let's look at the "classical" analysis situation. Say I have a sample of DST events (therefore these events have already been run through a production program) and I want to look for $t\bar{t}$ (TOP) candidates in the $e + \mu$ channel. I write a (Fortran) program, link it, and run it. Here I am assuming that these events are approximately 1 Mbyte long. The following is a play-by-play description of what happens.

1. Read in each event. If the events are on tape, access the tape first. For each event
 - Cut on global event quantities (e.g. primary vertex, trigger bit, etc.).
 - Require at least 1 e and 1 μ candidate.
 - Cut on e and μ quantities to get a "clean" sample.
 - Global event cuts (jets, E_t , etc.).
 - Plot p_{te} vs. $p_{t\mu}$.
2. If desired, save events passing e and μ p_t cuts onto disk or tape.

This procedure is very straightforward, and there are no obvious inefficiencies to be trimmed away. However, there are inefficiencies which involve the repetition of this procedure many times due to, for example, debugging the code or the application of additional cuts which result in smaller and smaller data sets. Events with a record size of 1 Mbyte (or more) can easily result in an unacceptable I/O load on the system, necessitating either a drastic reduction in the size of

each event through either a compression or trimming (miniDSTs, μ DSTs, etc.) and/or a drastic cut in the number of events.

An alternative to the “classical” way of storing data in event records consists of actually adding to the record length by introducing a structure to the data which would make it possible to read in only parts of the data record as needed. In other words, turn the data into a real “object-oriented” database. In the following sections, we present a layperson’s explanation of a relational and object-oriented database and how it may benefit the analysis of HEP data.

The Relational Database

The relational database was developed in the early 1970s, sparked by the influential paper of Codd[1]. Prototypes emerged in the middle to the late 1970s, including Ingres[2] and System R[3]. To visualize a **relational database**, consider a set of variables, or equivalently a multi-dimensional system space. In database parlance, each dimension (or variable) is referred to as a **domain**. Now consider that in this system, there are relationships, or conditions, between domains which tend to group certain values of each domain across domain boundaries. Such relationships group the domains into tables where the columns are domains and the rows (called **tuples**) are the grouping of domains. As an example, consider a HEP event at the high level where one is concerned with “objects” such as the 4-vectors of the electrons, muons, jets, and etc. One domain is the set of all electron transverse energies, another might be the same for muons. Each row is an event, which is the *relation* which causes particular electrons and muons to be grouped, and the entire set of relations forms a **table** (or **relation**)

Run	Event	Pte ($p_t e^+ / e^-$)	Ptmu ($p_t \mu^+ \mu^-$)
1001	126390	74.1	3.5
1006	69372	0	78.0
1007	574930	62.3	215.2
1020	6320	99.2	0
	⋮		

Figure 1. A “Typical” HEP Table Called **HEPdata**

in database parlance. In a relational database, the number columns is fixed. See Figure 1 for an example of a HEP table.

An important attribute of a relation is that it abstracts a certain type of file. A file of this type would be a sequence of records, one for each tuple in the database. Each record would consist of a sequence of fields, one for each column in the table. The definition of a relation implies that all records in the file would have the same number of fields, that no duplicate records be allowed, and that each field be atomic and have no additional structure.

Relational databases support special languages (called *query languages*) which provide for the analysis of the data in the database. For instance, in the HEP example, if you wanted to collect all events with an electron AND a muon with p_t above a cut (p_{cut}), the query language may take the form as in Figure 2 and produces a new table as in Figure 3. A query is simply a function on the database acting on relations and producing other relations. It is an essential feature of such queries that they produce other relations, which can then be

Range of t is HEPData
 Retrieve into TOPTable ($t.Run, t.Event$)
 Where $t.Pte > p_{cut}$ and $t.Ptmu > p_{cut}$

Figure 2. A Typical Query

Run	Event	Pte ($p_t e^+ / e^-$)	Ptmu ($p_t \mu^+ \mu^-$)
1007	574930	62.3	215.2
	⋮		

Figure 3. Result of Query on HEPdata

queried in exactly the same way. Relational databases also support the taking of intersections and unions of relations, which are called joins and meets, to produce other relations. In this way, a relational database provides a powerful means of completing queries on large databases containing numeric and alphabetic fields.

The Extensible Database

By the mid 1980s, it was recognized as a good idea to relax the requirement that domains consist of atomic objects and extend databases to support a collection of data types richer than merely numbers and strings. By the late 1980s, a number of such systems have been prototyped, including Iris [4], Orion [5], and GemStone [6]. Imagine a relational database which is extended to support not only numbers and characters but also physical (and logical) objects such as events,

muons, calorimetry tower lists, and etc. In other words, the table entries are no longer assumed to atomic, but rather to have an internal structure. For example, in Figure 1, imagine that there is a column (domain) called "TOP" which is a basic data type known to the system. Queries such as in Figure 4 would then be possible. It is important to note that:

- Queries always produce other tables, which can themselves be made the subjects of queries;
- When queries produce tables, the underlying data itself are not copied, but rather appropriate pointers are introduced pointing to the original objects, which may be quite large;
- Indices can be introduced to speed up subsequent queries.

The inclusion of new data types requires that new query operators be supported, new storage and access methods be developed, and new methods of optimizing queries be found. For example, it is important that the database support the storage and retrieval of events of an arbitrarily large size. It is also important that query operators defined on entire objects be supported, such as operators returning various statistical quantities. It is essential for HEP applications that arbitrary Fortran (and C) user-written subprograms interface easily into the query language, i.e. that there be a "hook" which the user can use to input customized operators to perform the queries. This is simply a logical extension of existing traditional HEP analysis methodology.

At this time, there is no consensus about how objects should be incorporated in a relational database. One viewpoint is to change the relational model as little

as possible and build a complete extensible database on top of it. For example, Stonebreaker [7] describes a mechanism for a user to register new abstract data types into the University version of Ingres. Because of the built-in access methods used by Ingres, the new data types must occupy a fixed amount of space. Still another viewpoint is to provide a modular and modifiable system on top of which extensible databases for specific applications can be built. The EXODUS database described by Carey [8] is an example of such an approach. See Carey [9] for a recent survey covering some of these issues and detailing other approaches.

Also note that although object-oriented programming is closely related to object-oriented databases, they are in fact two different concepts.

The Scientific Database

Beginning in the late 1970s and early 1980s, there has been increasing interest in statistical and scientific databases. These types of databases differ in a number of important ways from conventional relational databases.

1. The data types are different: data types in a statistical database include time series, matrices, and usually make provisions for missing data, outliers, etc.; data types in a conventional relational database are usually limited to numbers and strings.
2. Operators computing averages, standard deviations, regressions, etc. are important in a statistical database, while operators computing meets and joins are important in a conventional relational database.

3. Updating and modifying the data in a conventional relational database is frequent; on the other hand, in a statistical database, data may be added to the database, but is rarely changed.

Much of the data in a statistical or scientific database is static: updates are infrequent, but the queries are usually computationally very intensive. For example, the queries of interest typically cannot be answered from the stored relations and data; rather, some information typically needs to be computed, such as a regression. In other words, a typical query requires extracting features or some type of summary information from the database. In the past few years, there has been increasing interest in statistical and scientific databases. See, for example, Hammond [10] and Rafanelli [11] and the references cited therein.

A HEP Database

The analysis of HEP experiments would benefit from an object-oriented extensible database, supporting basic objects such as runs, events, candidates, etc. and allowing statistical and graphical queries.

The database would support *tuples* of objects, *relations* (expressing the relationships between tuples), and *operators* performed on the tuples and relations. The objects would include events, candidates, electrons, muons, jets, photons, etc. A relation could be viewed in many different ways: as a table, as a histogram, as a graph, or as an icon.

The operators would include:

Tabulate This would display the relation in a tabular format, not unlike a spreadsheet.

Graph This would graph the data, in any of a variety of formats.

Histogram This would provide various histograms associated with the tuples in a relation.

Iconize This would collapse the relation into an icon so that it could be moving around and so that other statistical operations could be performed on it.

Group This would allow the structure of the relations in a database, similar to the way files are structured into directories or prose is structured in an outline. That is collection of similar relations could be gathered together and collapsed into a single entity, until a later time that the underlying relations need to be accessed.

Mean This would compute the mean and other standard statistical functions of a relation.

Function fitting This would allow fitting functions to describe the functional relation of the data specified by tuples, regardless of the view of the data. That is, regardless of whether the data was displayed in a graphical, histogram, or iconic format.

Such a database should be developed with the following guidelines:

1. It should be designed using principles of object-oriented programming and coded in C++.
2. It would be developed in a UNIX environment and run under a graphics standard, such as X-windows.

3. It should be distributed, allowing for data to be distributed on a number of networked data servers, and queries to be processed on a number of network query servers.
4. It should allow for user-written (perhaps even Fortran?) “subroutines” to perform customized queries.

In a system with an object-oriented data structure, magnetic and/or optical tapes (see section 2.1.1) would be used for archiving purposes only, with all data and/or pointers to the data stored on random accessed devices such as optical disks. The advantage of disk over magnetic tape is in the ability to have a file structure. Access to files on disks which contain more than one file does not mean reading each preceding file. The organization of the physical data can be optimized to limit the amount of “unused” bits read in on each query (so-called query optimization). For instance, each domain can be structured like a file, and read in as needed. For instance, each domain can be structured like a file, and so only information (domains) which are needed for the query are read in. Another possibility is that frequently used collection of events (e.g. inclusive electrons, inclusive jets, etc.) can be processed into a table which contains pointers to a subset of all possible domains. Subsequent queries would be on these subsets of events. For **PRODUCTION** purposes, the entire collection of objects would be read into memory. Note that this software would not be hardware independent.

With the data in an object-oriented database, the above $t\bar{t} e\mu$ analysis is effectively reduced to a single query similar to Figure 2. This query would perform the following:

1. Retrieve for each event from a given file (table) information needed to calculate the number of electron and muon candidates.
2. Require at least one of these (passing cuts).
3. Save the list of events.

Once the list is saved, one could go back to the data file and read in global event "objects" for events in the list, make cuts, save the list of events passing cuts, and etc. Advantages of this scheme are:

1. The fact that one can now use the hardware and system software to modularize the analysis could significantly reduce the total amount of I/O.
2. A powerful bookkeeping system could be set up by adding a "bookkeeping" object during **PRODUCTION**. The object could be accessed independent of the rest of the event to keep track of event numbers, topology, etc.
3. Any subset of events could be copied onto disk from a list without sequential access of all events (since now an "event" is like a VAX "file" with a directory you can use to find it) and mailed to a collaborating institution.
4. With the help of a "PAW-like" analysis program which is written for this particular data structure, access to the data would be immediate and interactive, and analysis would consist of writing "command" files and "queries" to drive the "PAW-like" program.

Building a HEP Database

In the last section, we outlined some of the basic requirements of a HEP database. In this section, some of the issues involved in designing and building a HEP database will be discussed. The main issues will be the following:

1. A high performance HEP requires finding the right balance between primary (magnetic memory), secondary (magnetic and online optical disks); and tertiary (nearly online optical tapes and disks) storage devices. Performance can be improved by caching secondary storage in primary storage, and tertiary storage in secondary storage; and by using data indices residing in memory to speed up the retrieval of data. Although these are standard database techniques, they will have to be adapted and tuned to fit the needs of a HEP.
2. A high performance HEP requires balancing data transfer and CPU requirements of the database. To process data queries at acceptable levels requires queries be processed in a parallel or distributed fashion. Imagine that a 1 Tbyte database is divided into 100 buckets of 10 Gbytes each and that a \$10,000 RISC workstation was assigned to process queries for each bucket. Alternately, imagine that the 1 Tbyte database is divided into 10 buckets of 100 Gbytes each and that \$100,000 query processor is assigned to each. Most queries could be arranged so that each query processor could work independently, with only minimal communication required between them. Given these assumptions, a distributed query processor of the type just described would result in a nearly linear speed up in the queries. The hardware cost of such a system would be approximately \$1 million. It is essential that a HEP database be designed to

support distributed query processing as well as distributed access to the data.

3. Although there are many commercial databases available, it seems unlikely that any of them are suitable for a HEP. More likely, a HEP will incorporate various database tools and utilities as they become available.

In the remainder of this section, the last issue is discussed in more detail.

There are several ways of building a HEP database:

1. build it on top of a conventional commercial relational database,
2. build it on top of a commercial high performance relational database,
3. build it on top of an object-oriented commercial database,
4. build it using database tools and utilities, and
5. build a HEP from the bottom up.

There are several commercial databases available, including Sybase, Ingres, and Oracle. Option 1 would build the HEP on top of one of these databases. These databases do not support the basic data types and operators required of a HEP; moreover, they do not scale to allow for 1-10 Tbyte database imagined for the HEP. In general, they do not seem suitable for large, high performance scientific databases.

Option 3 would build the HEP on top of a commercial object-oriented extensible database. At present, such databases all allow for the creation of very general object types. This introduces performance penalties that are unacceptable for the HEP. Also, since these databases are just coming to the market now,

it is too early to commit to a company that may or may not be around six months from now, and to software that only represents a few man years of work.

Option 4 would build the HEP using tools and utilities designed to help shorten the cycle of designing and building a database. The main product of this type is Sun Microsystem's NETISAM. The idea would be to use whatever tools and utilities are available, but in the end to produce an extensible, object-oriented database that supports the data types and operators required for a HEP database.

Option 5 would design and build the HEP database from the bottom up. Such an approach would produce a database with the highest performance, but would make it more difficult to modify and maintain the code.

At present, it is too early to say with certainty which approach is the best, except that Option 1 is not viable. Our recommendation is to prototype Options 4 and 5 over the next one to two years and at the same time to monitor closely Options 2 and 3. It is possible that one of the latter two options could become viable within this timeframe.

3. REMOTE ACCESS: CLIENT/SERVER

It is probably universally agreed that the large SSC detector collaborations will want to have a computing group resident in some central location, and the traditional (and logical) choice will be the SSCL. In a scheme for computing as suggested above, access to the data by physicists resident at the SSCL would

be trivial. In the next section, we discuss the more non-trivial access, by non-resident collaborating members over network links.

3.1 COLLABORATING UNIVERSITY/LABORATORY ACCESS

Access to the data from outside the SSCL will present a real challenge to the implementation of any computing system. It will probably not be sufficient to require all data be accessed locally at SSCL while at the same time it will not be possible to furnish every collaborating institution with all of the data (certainly very little, if any, raw data will be distributed as a rule). It is envisioned that a reasonable compromise will be to allow both. Access to the entire data set at the SSCL should be made possible through high bandwidth, reliable links directly into the SSCL. Also, a scheme with “satellite” centers distributed uniformly (in bandwidth space) throughout the community (several per collaboration perhaps?) should be considered based on available resources. Access to subsets of the data (high level objects) should be via distribution of specific subsets of the data to collaborating institutions. For instance, if the optical disk is chosen as the primary storage medium, a second “juke-box” can be used to produce disk copies in batch at the SSCL for outside institutions. Access to this data will require one of the disk drives to be locally resident. The \simeq \$20,000 cost is not out of the grasp of even the smallest collaborating institution. However, the \simeq \$700 per disk might become a problem. One can imagine that the technology will catch up to this between now and the next eight years.

The above proposed scheme calls for turning the entire dataset into an object-oriented database. Software would have to be developed to do the queries on the

database, store any resultant "objects", and etc. Analysis of the data would be through a scheme much like the one used in the "PAW" program (as opposed to the traditional approach of writing, linking, and running code). The trend in private industry appears to be in the direction of distributed computing, using "clients" (local workstations) for interactive needs (e.g. make histograms and mathematical fits) and "servers" (host computers with access to the data) sending data over high-speed links. Such a scheme, although attractive in many ways, depends *entirely* on the reliability of the data links. For this reason, consideration should be given to providing sufficient redundancy through the "satellite" centers (see above). Note that the client/server approach has already been standardized in the MIT program "X-WINDOWS" system, and MICROSOFT has recently committed itself to this direction.

4. SSCL IMPLICATIONS

4.1 SSCL COMPUTER SYSTEM

Two of the largest HEP computing centers in the U.S., SLAC and FNAL, are interesting examples of two very different ways to solve the same problem. SLAC relies on powerful, well-maintained and supported mainframes available to all experimenters. At FNAL, on the other hand, much (although certainly not all) of the computing is done at collaborating institutions, an informal and primitive "distributed computing" arrangement involving many different machine architectures and operating systems. One of the more pressing issues the new SSCL faces concerns the use of the central-mainframe(s) versus distributed computing

(workstations). In the above scheme, however, both will have to be used. One needs a large central computer center with sufficient I/O bandwidth organized around the large amounts of data. These computer(s) would act as servers to a farm of clients for **PRODUCTION** of DSTs and for CPU intensive simulation work as well as servers to any client (workstation) at a collaborating institution for remote access and analysis. [Note that such client (workstations) will invariably have enough CPU power to handle most analysis requirements if one assumes that a sub-\$10,000 workstation in the 10-20 MIP range would be readily available.] The proper segmentation of such a computer system around I/O, serving, and CPU-intensive tasks should be studied. One would suspect that there is a maximum efficiency for some form of segmentation, providing that the specifications of the computer system are well established.

4.2 EXISTING RESOURCES

Another, perhaps more political, issue concerns the problem of what to do with the current tremendous investment in and variety of computing at the universities and national laboratories. A "PAW"-like system would help alleviate the issue. The idea of "operating system" may very well, at some point in the future, disappear.

4.3 KEY FACTORS

The following is a list of key areas which are needed for the above computing scheme to be successful:

1. *Networks.* The proliferation of T1 (or greater) cable links is good news. It will be extremely important that there be easy access to this network by any collaborating institution, and that the network be reliable. This may very well be the weakest link in principle, however it is clear that the importance of powerful networks has not been lost on the computer industry.
2. *Data Storage.* No more tapes, except maybe for archival storage. Serious resources (money) should be spent on file-oriented disks with the goal that at least 1 Tbyte of data be accessible interactively. Projected technological advances makes this a reasonable goal.
3. *Software.* There is a large amount of code to write to set up such a system but the basics are already there (X-windows from MIT, Windows at Microsoft, etc.). An interactive "PAW"-like system would have the advantage that if it was supported in UNIX, VMS, VM, and MS/DOS, there would probably be no other operating system one would have to consider. (In eight years the list could be even shorter than this.) The client/server scheme is reported to be operating system independent to some degree. The "PAW"-like system would necessitate creating the following:
 - (a) A "PAW"-like shell. In principle, there is no reason that the current PAW program could not be modified for this.
 - (b) An extensible relational database to store, retrieve, and access the data, as described in section 2.3.1.

- (c) A database query language to connect "PAW" to the data. The design of an extensible, relational database together with the appropriate query language is a large task. There are presently several DOE and NASA funded software projects with similar aims, which could either be extended or absorbed if appropriate.

Data taking is scheduled to begin sometime during 1998. The detector builders claim that it will take at least this long to build the detector. In implementing the above scheme, one would not have eight years of lead time since software systems will be needed as early as possible to set up a working collaboration. A powerful, modern system as just outlined could be functioning several years before data dating if enough resources are provided. And such a system could, in principle, be utilized by all of the collaborations which will do physics at the SSCL, obviating the need for duplicating the work of setting up a different **ONLINE, PRODUCTION, and ANALYSIS** system for each detector. An additional justification can be made that such a system would be extremely powerful and available to both the scientific and business community—a beneficial HEP spinoff.

REFERENCES

1. E. Codd, *A relational model of data for large shared data banks*, Com. ACM, **13-6**, pp. 377-387.
2. *The Ingres Papers*, M. Stonebreaker, ed., Addison-Wesley, 1986.
3. Astrahan et al., "System R: A Relational Approach to Database Management," *ACM Transactions on Database Systems*, **1** (1976).
4. D. Fishman et al., "Overview of the Iris DBMS," *Object-Oriented Concepts, Databases, and Applications*, W. Kim and F. H. Lochovsky, eds., ACM, New York, 1989, pp. 219-250.
5. W. Kim et al., "Features of the ORION Object-Oriented Database System," *Object-Oriented Concepts, Databases, and Applications*, W. Kim and F. H. Lochovsky, eds., ACM, New York, 1989, pp. 251-282.
6. R. Bertl et al., "The GemStone Data Management System," *Object-Oriented Concepts, Databases, and Applications*, W. Kim and F. H. Lochovsky, eds., ACM, New York, 1989, pp. 283-308.
7. M. Stonebreaker, "Inclusion of New Types in Relational Data Base Systems," *Proceedings of IEEE/Data Engineering*, IEEE, 1986, pp. 262-269.
8. M. J. Carey, "The Architecture of the EXODUS Extensible DBMS," *Proceedings of the Object-Oriented Database Workshop*, ACM, 1986, pp. 52-65.
9. M. J. Carey, "Special Issue on Extensible Database Systems," *Database Engineering*, 1987.

10. *Proceedings of the Second International Workshop on Statistical Database Management*, R. Hammond and J. L. McCarthy, eds., Springer, 1983.
11. "Research Topics in Statistical and Scientific Database Management," in *Statistical and Scientific Database Management*, M. Rafanelli, J. C. Klensin, and P. Svensson, eds., Springer, Berlin, 1989.