# Lawrence Berkeley National Laboratory

**Title**
Design of a Very Large Storage System

**Permalink**

**Authors**
Penny, Samuel J
Fink, Robert
Alston-Garnjost, Margaret

**Publication Date**

# DESIGN OF A VERY LARGE STORAGE SYSTEM

Samuel J. Penny, Robert Fink, and
Margaret Alston-Garnjost

June 1970

LAWRENCE RADIATION LABORATORY
UNIVERSITY of CALIFORNIA BERKELEY

UCRL-19757 Rev.

# DISCLAIMER

# DESIGN OF A VERY LARGE STORAGE SYSTEM

by Samuel J. Penny, Robert Fink, and Margaret Alston-Garnjost

Lawrence Radiation Laboratory, University of California

Berkeley, California

The Mass Storage System (MSS), built around the IBM 1360 photo-digital store, offers the users of the CDC 6600 computer complex at the Lawrence Radiation Laboratory—Berkeley, a facility for the storage, management, and retrieval of very large volumes of data. The system is capable of storing $3 \times 10^{11}$ bits of data on-line to the computer. Maximum access time is 5 seconds. Off-line archival storage is unlimited and data can be moved back on-line within a short time. An automated data-management system designed to handle these large amounts of data offers many new benefits and possibilities. Reliable on-line data bases of a very large size open the way to new and interesting applications as well as new ways to approach existing problems.

(word count 107)

# DESIGN OF A VERY LARGE STORAGE SYSTEM[*]        Paper No. 087

by Samuel J. Penny, Robert Fink, and Margaret Alston-Garnjost

Lawrence Radiation Laboratory, University of California

Berkeley, California 94720

## INTRODUCTION

The Mass Storage System (MSS) is a data-management system for the on-line storage and retrieval of very large amounts of permanent data. The MSS uses an IBM 1360 photo-digital storage system (called the chipstore) with an on-line capacity of $3 \times 10^{11}$ bits as its data storage and retrieval equipment. It also uses a CDC 854 disk pack for the storage of control tables and indices. Both these devices are attached to a CDC 6600 digital computer at the Lawrence Radiation Laboratory—Berkeley.

---

[*]Work done under auspices of the U. S. Atomic Energy Commission

---

Plans for the MSS began in 1963 with a search for an alternative to magnetic tape as data storage for analyses in the field of high energy physics. A contract was signed with IBM in 1965 for the chipstore, and it was delivered in March of 1968. The associated software on the 6600 was designed, produced, and tested by LRL personnel, and the Mass Storage System was made available as a production facility in July of 1969.

This paper is concerned with the design effort that was made in developing the Mass Storage System. The important design decisions,

and some of the reasons behind those decisions, are discussed.  Brief
descriptions of the hardware and software illustrate the final result of
this effort.

## CHOICE OF THE HARDWARE

By 1963 the analysis of nuclear particle interactions had become
a very large application on the digital computers at LRL–Berkeley.
More than half the available time on the IBM 7094 computer was being
used for this analysis, and the effort was expanding.  Much of the
problem was purely data manipulation -- sorting, merging, scanning,
and indexing large tape files -- and single experiments produced tape
libraries of hundreds of reels each.

The problems of handling large tape libraries had become well
known to the experimenters.  Tapes were lost; they developed bad spots;
the wrong tapes were used; keeping track of what data were on what tape
became a major effort.  All these problems degraded the quality of the
data and made the experiments more expensive.  A definite need existed
for a new approach.

The study of the problem began with establishment of a set of
criteria for a large-capacity on-line storage device, and members of
the LRL staff started investigating commerically available equipment.
The basic criteria used were:

(a)  The storage device should be on-line to the central computing
facility.

(b) It should have an on-line capacity of at least $2.5 \times 10^{11}$ bits (equivalent to 2000 reels of tape).

(c) Access time to data in the storage device should be no more than a few seconds.

(d) The data-reading transfer rate should be at least as fast as magnetic tape.

(e) The device should have random-access capability.

(f) The storage medium of the device should be of archival quality, lasting 5 years at least.

(g) The storage medium need not be rewritable.

(h) The frequency of unrecoverable read errors should be much lower than on magnetic tape.

(i) Data should be easily movable between the on-line storage device and shelf storage.

(j) The device hardware should be reliable and not subject to excessive failures and down time.

(k) Finally, the storage device should be economically worthwhile and within our budget.

Several devices were proposed to the Laboratory by various vendors. After careful study, including computer simulation of the hardware and scientific evaluations of the technologies, the decision was made to enter into a contract with IBM for delivery, in fiscal year 1968, of the 1360 photo-digital storage system. This contract was signed in June of 1965.

The major application contemplated at that time is described in Ref. 1.

It was clear that one of the major problems in the design of the associated software would be the storage and maintenance of control tables and indices to the data. Unless indexing was handled automatically by the software, the storage system would quickly become more of a problem than it was worth. Protection of the indices was seen to be equally important, for the system would be dependent on them to physically locate the data. It was decided that a magnetic disk pack drive, with its removable pack, was the most suitable device for the storage of the MSS tables and indices.

A CDC 854 disk pack drive was purchased for this purpose.

## DESCRIPTION OF THE HARDWARE

### 1360 Photo-digital storage system

The IBM 1360 chipstore is an input-output device composed of a storage file containing 2250 boxes of silver halide film chips, a chip recorder-developer, and a chip reader. Figure 1 shows the general arrangement of the chipstore hardware and its relation to the CDC 6600 computer. References 2 through 5 describe the hardware in detail. A brief summary is given below.

A chip is 35 by 70 mm in size and holds 4.7 million bits of data as well as addressing and error-correction or error-detection codes. Data from the 6600 computer are recorded on the chip in a vacuum with an

electron beam, taking about 18 sec per chip. The automatic film devel-
oper unit completes the processing of a chip within 2.5 min; it overlaps
the developing of eight chips so that its processing rate is comparable
to that of the recorder.

Up to 32 chips are stored together in a plastic box. Figure 2 shows
a recorded film chip and the box in which it is kept. These boxes are
transported between the recorder-developer, the box storage file, and
the chip reader station by means of an air blower system. Transport
times between modules on the Berkeley system average around 3 sec.

Under the command of the 6600 computer the chipstore transports
a box from the storage file to the reader, picks out a chip, and positions
it for reading. The chip is read with a spot of light generated by a
cathode-ray tube and detected by a photomultiplier tube at an effective
data rate of 2 million bits per second. The error correction-detection
codes are checked for validity as the data are read, and if the data are
incorrect, an extensive reread and error-correction scheme is used to
try to reproduce the correct data. The data are then sent to the 6600
across a high-speed data channel. Chip pick and store times are less
than 0.5 sec.

The box storage file on the Berkeley 1360 system has a capacity
of 2250 boxes. This represents an on-line data capacity of 2750 full
reels of magnetic tape (at 800 BPI); 1360 systems at other sites have
additional file modules, giving them an on-line capacity three or more

times as great as at Berkeley.

A manual entry station on the chipstore allows boxes of chips to be taken out of the system or to be reinserted.  By keeping the currently unused data in off-line storage and retaining only the active data in the file, the potential size of the data base that can be built in the MSS is equivalent to tens of thousands of magnetic tapes.

A process control computer is built into the chipstore hardware. This small computer is responsible for controlling all hardware actions as well as diagnosing malfunctions.  It also does the detailed scheduling of events on the device.  Communication between the chipstore and the host computer goes through this processor.  This relieves the host of the responsibility of commanding the hardware in detail, and offers a great deal of flexibility.

## 854 Disk pack drive

The CDC 854 disk pack drive holds a removable 10-surface disk pack.  The pack has a typical access time of 90 msec, and a data transfer rate of about 1 million bits per sec.  Its storage capacity is 48 million bits.

MSS uses this pack for the storage of all its tables and indices to the data that have been written into the 1360 chipstore.  A disk pack was chosen for this function to insure the integrity of the MSS tables.  The 854 has a proven record of hardware and data reliability.  Also, since the pack is removable, the drive can be repaired and serviced without threat to the tables.

## 6600 Computer complex

The chipstore is connected to one of the CDC 6600 computers at LRL through a high-speed data channel. The 6600 computer has 131 072 words of 60-bit central core memory (CM), a central processor unit (CPU) operating at a 100-nsec cycle rate, and 10 peripheral processor units (PPU). Each PPU contains 4096 words of 12-bit core memory and operates at a 1-msec cycle rate. The PPU's control the data channel connections to the external input-output equipment, and act as the interface between jobs residing in CM and the external world.

The operating system on the 6600 is multiprogrammed to allow several jobs to reside in CM at once and share the use of the CPU. Two of the PPU's act as the system monitor and operator interface for the system, and those remaining are available to process task requests from the monitor and execute jobs. The MSS, composed of both CPU and PPU programs, has been built as a subsystem to this operating system.

## CHOICE OF THE MASS STORAGE SYSTEM SOFTWARE

## Design objectives

Having made the commitment on hardware, the Laboratory was faced with designing and implementing the associated software. The basic problem was to produce a software system on the CDC 6600 computer that, using the IBM 1360 chipstore, would lead to the greatest increase in the productive capacity of scientists at the Laboratory. In

addition, it was necessary that the system be one that the scientists

would accept and use, and to which they would be willing to entrust their

data. It would be required to be of modular design and "open-ended,"

allowing expansion and adjustment to new techniques that the scientists

might develop for their data analysis.

Overall study of the problem yielded three primary objectives.

Most important was to increase the reliability of the data storage, both

by reducing the number of data-read errors and by protecting the data

from being lost or destroyed; much time and effort could be saved if this

objective were met. The second objective was to increase the utilization

of the whole computer complex. The third was to provide facilities for

new, more efficient approaches to data analysis in the future.

The problem was divided into three technical design areas: the

interaction between the software and the hardware, the interaction

between the user and the software, and the structure of the stored data.

In the area of software-hardware interaction, the design objectives

were to maximize protection of the user data, interleave the actions for

several jobs on the hardware, reduce the need for operator intervention,

and realize maximum utilization of the hardware. This was the approxi-

mate order of importance.

Objectives in the area of user interaction with the MSS included

making that interaction easy for the user, offering him a flexible data-

read capability, and supplying him with a protected environment for

his data. Ease of data manipulation was of high value, but not at the
expense of data protection. A flexible read mechanism was necessary,
since if the users could not read their data from the MSS, they would
seek other devices. This flexibility was to include reading data from
the chipstore at rates up to its hardware limit, having random access
to the data under user control, possibly intermixing data from the chip-
store, magnetic tapes, and system disk files, and being able to read
volumes of data ranging in size from a single word to the equivalent of
many reels of tape.

The problem of data structures for the MSS was primarily one of
finding a framework into which existing data could be formatted and which
met the requirements of system and user interaction. This included the
ability to handle variable-length data records and files and to access
these data in a random fashion. It was decided that a provision to let
the user reference his data by name and to let the system dynamically
allocate storage space was very important. It was also important to
have flexible on-line—off-line data-transfer facility so that inactive data
could be moved out of the way.

## Software design decisions

Several important design decisions were made that have had a
strong effect on the nature of the final system. Some of these decisions
are listed here.

Each box used for data storage is given a unique identification number, and this number appears on a label attached to the box. A film chip containing data is given a unique home address, consisting of the identification number of the box in which it is to reside and the slot in that box where it is to be kept. Control words written at the beginning of the chip and at various places throughout the data contain this address (along with the location of the control word on the chip), and this information can be checked by the system to guarantee correct positioning for retrieval of the data. It is also used to aid in recovery procedures for identifying boxes and chips. This control information can be used to help reconstruct the MSS tables if they are destroyed.

The control words are written in context with the data to define the record and file structure of the data on the chips. The user is allowed to give the address of any control word (such as the one at the beginning of a record) to specify what data are to be read. This scheme meets the design objective of allowing random access to data in the chipstore.

Data to be written into the chipstore are effectively staged. The user must have prepared the data he wishes to be recorded in the record and file structure he desires in some prior operation. He then initiates the execution of a system function that puts the source data into chip format, causes its recording on film chips, waits for the chips to be developed, does a read check of the data, and then updates the MSS tables.

Data read from the chipstore are normally sent directly to the user's program, though system utility functions are provided for copying data from the chipstore to tape or disk. If the user desires, he may include a system read subroutine with his object program that will take data directly from the chipstore and supply them to his executing program. This method was chosen to meet the objectives of high data-transfer rates and to provide the ability to read gigantic files of data.

To aid the user in the access and management of his data in the MSS, it was decided to create a data-management control language oriented to applications on the chipstore. A user can label his data with names of his own choosing and reference the data by those names. A two-level hierarchy of identification is used, that of data set and subset. The data set is a collection of named subsets, in which each subset is some structure of user data. The control language is not limited to manipulating only data from the chipstore; it can also be used to work with magnetic tape or system disk files.

Two more decisions have greatly simplified the overall problem of data management in the MSS. The first was to allocate most of the on-line storage space on the chipstore in blocks to the scientists engaged in data analysis of current experiments, and give them the responsibility of choosing which of their data are to reside on-line within their block and which are to be moved off-line. The second decision was to treat all as permanent. Once successfully written, film chips are never physically destroyed. At most, the user may delete his reference to the data, and the chips are moved off-line.

DESCRIPTION OF THE MSS SOFTWARE

The system in use on the 6600 computer for utilizing the chipstore results both from design effort at the beginning of the project and from experience gained during the implementation and initial production phases. Its essential features are listed below.

Indexing and control of the data stored in the chipstore are handled through five tables kept on the disk pack, as follows.

The box group allocation table controls the allocation of on-line storage space to the various scientists or experiments at the Laboratory. Any attempt by a user to expand the amount of on-line space in use by his box group above its allowable limit will cause his job to be aborted.

The box identification table contains an entry for each uniquely numbered box containing user data chips. An entry tells which box group owns the box, where that box is stored (on-line or off-line), which chip slots are used in the box, and the date of its last use.

The file position table describes the current contents of the 1360 file module, defines the use of each pocket in the file, and gives the identification number of the box stored in it.

The data set table contains an entry for each of the named collections of data stored in the chipstore. Status and accounting information is kept with each data set table entry. Each active entry also points to the list of subsets collected under that data set.

The subset list table contains the lists of named subsets belonging to the entries in the data set table. A subset entry in a list gives the name

of the subset, the address of the data making up that subset, and status information about the subset.

These tables are accessed through a special PPU task processor program called DPR. This processor reads or writes the entries in the tables as directed. However, if the tables are to be written, special checks and procedures are used to aid in their protection. Twice daily the entire contents of the MSS disk pack are copied onto magnetic tape. This is backup in case the data on the pack are lost.

All communication to the chipstore across the data channel link is handled through another PPU task processor program called 1CS; 1CS is multiprogrammed so that it can be servicing more than one job at a time. Part of its responsibility is to schedule the requests of the various user jobs to make most effective use of the system. For instance, jobs requiring a small amount of data are allowed to interrupt long read jobs. Algorithms for overlapping box moving, chip reading, and chip writing are also used to make more effective use of the hardware.

1CS and DPR act as task processors for jobs residing in the central memory of the 6600. The jobs use the MSSREAD subroutine (to read from the chipstore) or the COPYMSS system utility to interface to these task processors. These central memory codes are described below.

The reading of data from the chipstore to a job in central memory is handled by a system subroutine called MSSREAD. The addresses of the data to be read and how the data are to be transmitted are given to

MSSREAD in a data-definition file. This file is prepared prior to the
use of MSSREAD by the COPYMSS program described later. MSSREAD
handles the reading of data from magnetic tape, from disk files, or from
the chipstore. If the data address is the name of a tape or disk file,
MSSREAD requests a PPU to perform the input of the data from the device
a record at a time. If the address is for data recorded in the chipstore,
it connects to 1CS, and working with that PPU code, takes data from the
chipstore, decodes the in-context structure, and supplies the data to the
calling program.

A system program called COPYMSS is responsible for supplying
the user with four of the more common functions in MSS. It processes
the MSS data-management control language to construct the data-definition
file for MSSREAD. It performs simple operations of copying data from
the chipstore to tape or disk files. It prepares reports for a user,
listing the status of his data sets and subsets. Finally, COPYMSS is
the program that writes the data onto film chips in the chipstore.

To write data to the chipstore, the user must prepare his data in
the record and file structure he desires. He then uses the MSS control
language to tell COPYMSS what the data set and subset names of the data
are to be and where the data can be found. COPYMSS inserts the required
control words as the data are sent through 1CS to the chipstore to be
recorded on film chips. After the chips have been developed, 1CS
rereads the data to verify that each chip is good. If a chip is not

recorded properly, it is discarded and the same data are written onto

a new chip. When all data have been successfully recorded and the chips

are stored in the home positions, COPYMSS uses DPR to update the disk

pack tables, noting the existence of the new data set—subset.

The remaining parts of the MSS software include accounting pro-

cedures, recovery programs, and programs to control the transfer of

data between on-line and off-line storage. These programs, used by

the computer operations group, are not available to the general user.

## RESULTS AND CONCLUSIONS

### Effort

A total of about 7.5 man-years of work was invested in the Mass

Storage System at LRL—Berkeley. The staff on the project was composed

of the authors with some help from other programmers in the Mathe-

matics and Computing Department. The breakdown of this effort is shown

in Table I.

Table I. Distribution of MSS implementation effort.

| Operation | Man-years |
|---|---|
| Procurement and Evaluation | 1.0 |
| System design | 2.8 |
| Software coding | 1.7 |
| Software checkout | 0.8 |
| Maintenance, documentation, etc. | 1.2 |

## Operating experience

The Mass Storage System has been in production status since June 1969. Initial reaction of most of the users was guarded, and many potential users were slow in converting to its use. As a result, usage was only about 2 hours a day for the first 3 months. Soon after, this level started to increase, and at the end of one year of production usage a typical week (in the month of June 1970) showed the usage given in Table II.

|  Table II.   MSS usage per week. | |
| --- | --- |
| Number of read jobs | 250 |
| Number of write jobs | 100 |
| Chips read | 11 500 |
| Bits read | $5.4 \times 10^{10}$ |
| Unrecoverable read errors | 15 |
| Chips written | 1 900 |
| Percentage down time | 8.5 |

Most of the reading from the chipstore is of a serial nature, though the use of the random-access capability is increasing. Proportionally more random access activity is expected in the future as users become more aware of its possibilities.

A comparison of the MSS with other data-storage systems at the Laboratory, shown in Table III, points out the reasons for the increased

Table III.  Comparison of storage devices at LRL–Berkeley.

| | MSS | CDC 607 tape drive | CDC 854 disk pack | IBM 2311 data cell | CDC 6603 system disk |
|---|---|---|---|---|---|
| On-line capacity (bits/device) | $3.3 \times 10^{11}$ | $1.2 \times 10^8$ | $4.8 \times 10^7$ | $3.0 \times 10^9$ | $4.5 \times 10^8$ |
| Equivalent reels of tape | 2750 | 1 | 0.4 | 25 | 3.75 |
| Cost of removable unit | $13/box | $20/reel | $500/pack | $500/cell | --- |
| Storage medium cost ($¢/10^3$ bits) | 0.008 | 0.017 | 1.0 | 0.17 | --- |
| Average random access time (sec) | 3 | (minutes) | 0.075 | 0.6 | 0.125 |
| Maximum transfer rate (kilobits/sec) | 2000 | 720 | 1330 | 450 | 3750 |
| Effective transfer rate[a] | 1100 | 500 | --- | 200 | 400 |
| Approximate capital costs (thousands of dollars) | 1000 | 100 | 35 | 220 | 220 |
| Mean error-free burst length (bits) | $1.6 \times 10^9$ | $2.5 \times 10^7$ | $> 10^{10}$ | $10^9$ | $> 10^{10}$ |

a. Based on usage at LRL–Berkeley; the rates given include device-positioning time.

usage. For large volumes of data, the closest competitor is magnetic tape (assumed here to be full 2400-foot reels, seven-track, recorded at 800 BPI).

The values shown in Table III are based on the following assumptions: on-line capacities are based on having a single unit (e. g. , a single tape drive); capital costs are not included in the storage medium costs; effective transfer rates are based on usage at LRL, and are very low for the system disk because all jobs are competing for its use; and all costs given are only approximate.

The average data-transfer rate on long read jobs (involving many chips and many boxes) is more than one million bits per second. This is decidedly better than magnetic tape. Short reads go much faster than from tape once the 3-sec access time is complete.

The biggest selling point for the Mass Storage System has been the extremely low data-error rate on reads. This rate is less than 1/60 of the error rate on magnetic tape. The second most important point has been the potential size of the data files stored in the chipstore. Several data bases of from 20 to 200 boxes of data have been constructed. Users find that having all their data on-line to the computer and not having to rely on the operators to hang tapes is a great advantage. Their jobs run faster and there is less chance that they will not run correctly.

The cost of storing data on the chipstore has proven to be competitive with magnetic tape, especially for short files or for files that will be

read a number of times. Users are beginning to find it profitable to
store their high-use temporary files on the chipstore.

The system has not been without its difficulties. Hardware reli-
ability has at times been an agonizing problem, but as usage increases
and the engineers gain more experience on the hardware, the down time
for the system has decreased significantly. We now feel that 5% down
time would be acceptable, though less would be preferable. Fortunately,
lack of hardware reliability has not affected the data reliability.

## CONCLUSIONS

Though intended primarily as a replacement for magnetic tape in
certain applications, the MSS has shown other benefits and capabilities.
Data reliability is many times better than for magnetic tape. Some
applications requiring error-free storage of large amounts of data simply
are not practical with magnetic tape, but they become practical on the
chipstore. The nominal read rate is faster than that of magnetic tape
for long serial files. In addition, any portion of a file is randomly
accessible in a time ranging from a few milliseconds to 5 seconds.

The MSS is not without its limitations and problems. The 1360 is
a limited-production device: only five have been built. It uses tech-
nologies within the state of the art but not thoroughly tested by long
experience. Keeping the system down time below reasonable limits is
a continuing and exacting effort. Development of both hardware and
software has been expensive. The software was a problem because the

chipstore was a new device and people had no experience with such large

storage systems.

The Mass Storage System has met its purpose of increasing the

productive capacity of scientists at the Laboratory. It has also brought

with it a new set of problems, as well as a new set of possibilities. The

biggest problem is how to live with a system of such large capacity, for

as more and more data are entrusted to the chipstore, the potential loss

in case of total failure increases rapidly. The MSS offers its users

important facilities not previously available to them. More important,

the age of the very large Mass Store has been entered. In the future, the

MSS will become an important tool in the computing industry.

## REFERENCES

1)    MARGARET H. ALSTON and SAMUEL J. PENNY

      The use of a large photodigital mass store for bubble chamber
      analysis
      IEEE Trans. Nucl. Sci. Volume NS-12 [4], Pages 160-163, 1965.

2)    J. D. KUEHLER and H. R. KERBY

      A photo-digital mass storage system
      AFIPS Conference Proceedings of the Fall Joint Computer Con-
      ference, Volume 29, Pages 735-742, 1966.

3)    I. B. OLDHAM, R. T. CHIEN, and D. T. TANG

      Error detection and correction in a photo-digital storage system
      IBM J. Res. Develop. Volume 12 [6], Pages 422-430, 1968.

4)    D. P. GUSTLIN and D. D. PRENTICE

Dynamic recovery techniques quarantee system reliability
AFIPS Conference Proceedings of the Fall Joint Computer Con-
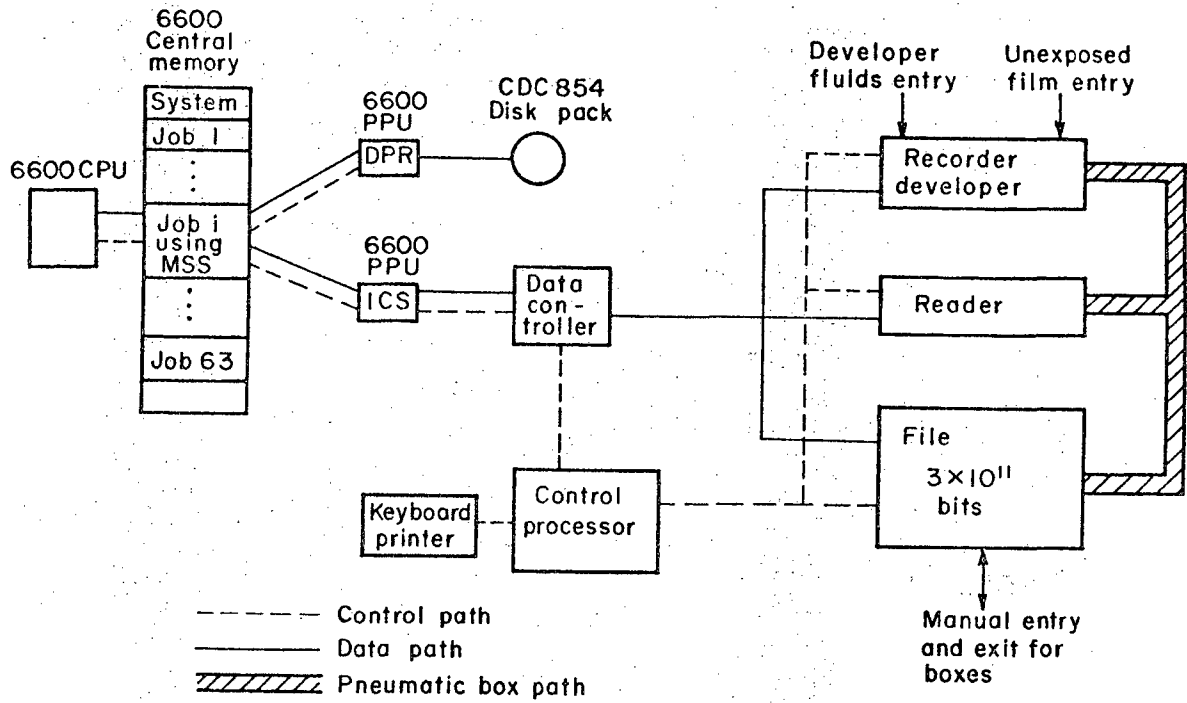ference, Part II, Volume 33, Pages 1389-1397, 1968.

5)    RICHARD M. FURMAN

IBM 1360 photo-digital storage system
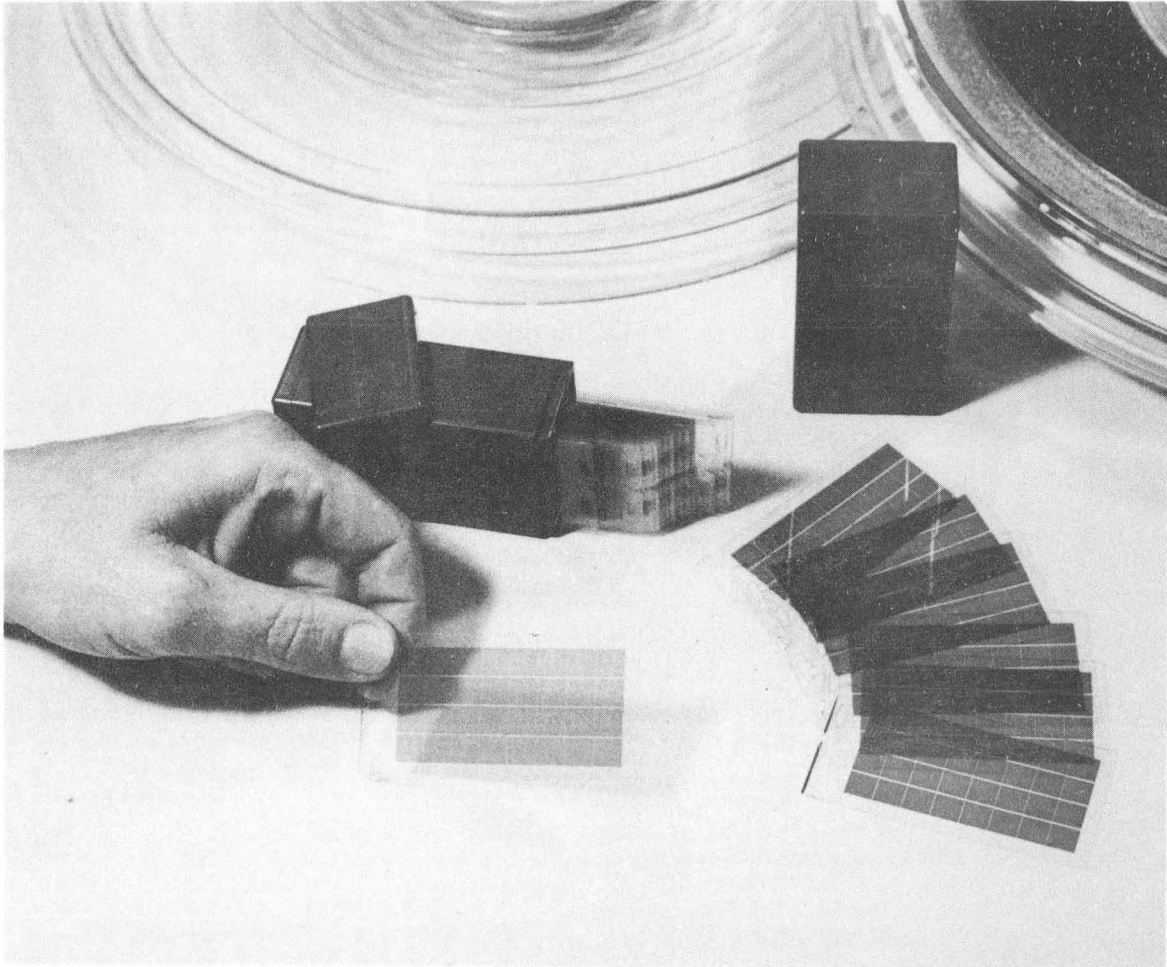IBM Technical Report TR 02.427, May 15, 1968.

FIGURE CAPTIONS

Figure 1.  General MSS architecture.

Figure 2.  Recorded film chips and storage box.

Fig. 1

XBB 689-5629

Figure 2.

TECHNICAL INFORMATION DIVISION
LAWRENCE RADIATION LABORATORY
UNIVERSITY OF CALIFORNIA
BERKELEY, CALIFORNIA 94720