

N. Armenise and A. Silvestri: DATA MANAGEMENT IN HIGH STATISTIC EXPERIMENTS. -

1. - INTRODUCTION. -

The actual effort in the experimental field of elementary particles is to make available experiments with a very large number of events ($\sim 10^6$).

Automatic devices developed to measure bubble-chamber events (HPD, PEPR, Spiral reader) and experimental devices of very high efficiency (Missing mass spectrometer, OMEGA project, etc.) seem very able to furnish experimental data at a high production rate. On the contrary, the traditional tools of management and analysis of the information in this field (experimental data) are critical points in the chain; they slow down the rate of data elaboration, and therefore limit the efficiency obtained in the first step.

Therefore a different organization of the second step is needed; this paper outlines the problem of automatic management of experimental data and presents a few realizations from H. E. Group of Bari University.

2. - THE FUNCTION OF SYSTEM'S MANAGEMENT. -

The different steps of management from the experimental phase (the exposure of an optical chamber to a beam of particles, or data collection by counters) to the DST (Data Summary Tape) can be "grosso modo" listed in the following way:

- a) Scan - I. e. a visual search for integral devices (see B. C.), a trigger for differential devices (such a spectrometer, counters, and so on); in this case, this step practically coincides with next step;
- b) Measurement;
- c) Geometrical and kinematical reconstruction;

2.

d) Selection of events and classification of sub-systems.

When the DST is made, the step of analysis begins. Fig. 1 displays the flux-diagram of these steps.

3. - CRITICAL POINTS. -

In the flux diagram of Fig. 1, a number in a circle identifies the critical steps in a "by hand" management of events. In fact the managements of large number of events present great difficulties for:

a) The identification of events badly measured because of topology or by accident. This identification is made with the aim to remeasure these events. This is feasible at least for bubble-chamber experiments. Alternatively the cut of bad events can lead to systematical biases, especially for the events rejected because of difficult topology;

b) The identification of events with failures in reconstruction. Generally the failure of events, too, is due to errors in the measurement;

c) The classification of successfull events and collection of the informations, useful in order to characterize a physical hypothesis without ambiguity;

d) The division of misidentified events in different subsystems.

Points c) and d), that require the management of a large number of events, are necessary to define a good sample (without systematical biases) and to give correct cross sections.

It is noticeable that a semiautomatic management, using sequential supports to record the useful information for the critical points, is able to handle efficiently only a small sample of events.

For instance the classification of events do require other informations than those available in the reconstruction. This fact generally implies the comparison between sequential files (two or more), which hold a few thousands data; consequently the time for this comparison becomes large.

4. - A FEASIBLE METHOD OF AUTOMATICAL ANALYSIS. -

Fig. 2 shows the flux diagram of the processing system of events actually used in H. E. group of Bari University. The critical points are handled automatically as far as this can be allowed by our non automatic measuring machines.

At the end of sequential flux, cross sections and other quantities of physical interest can be obtained without ambiguity and errors.

The method used, is essentially based on a direct access file, in which each event corresponds to a record addressable by means of a Key

DIAGRAM FLUX OF STEPS CONCERNING TO EVENTS MANAGEMENT

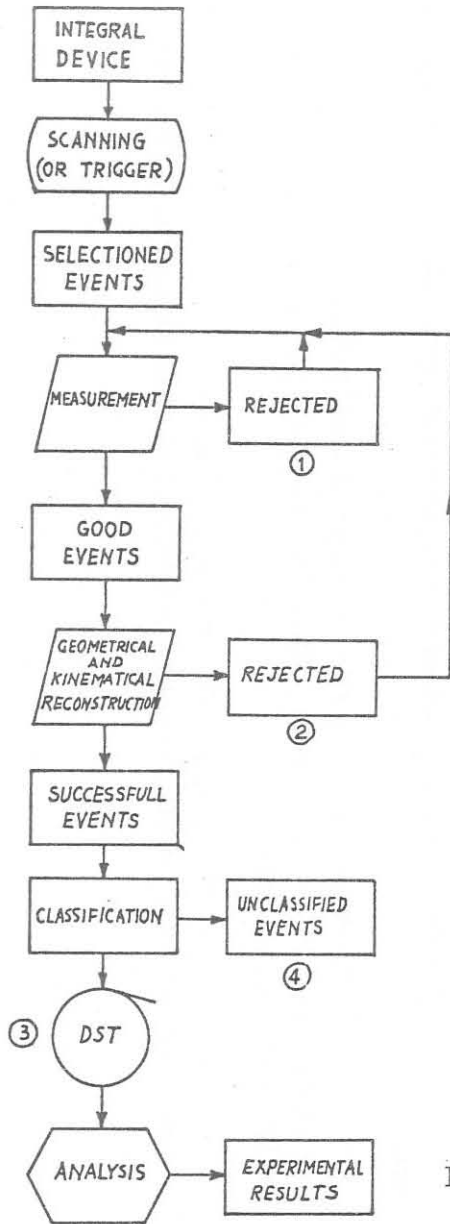


FIG. 1

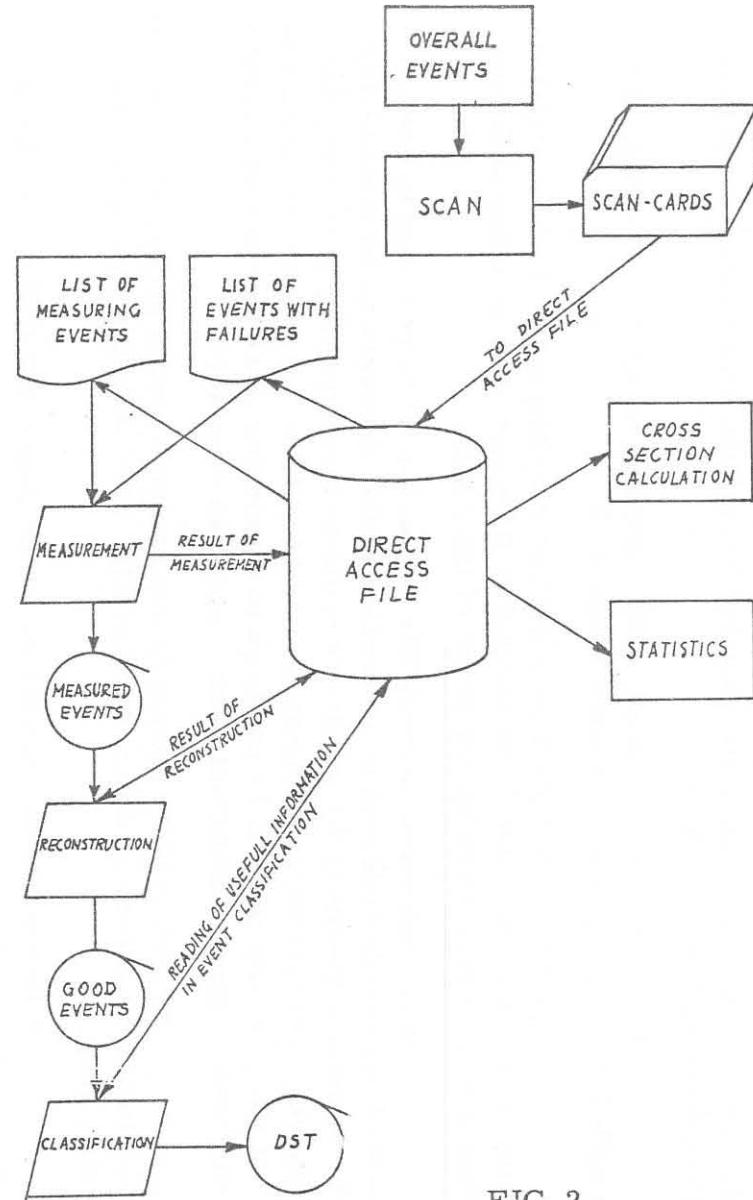


FIG. 2

4.

(Number of roll, Frame of event, the position of event in the photo). The record holds informations upon the topology of events (number of tracks, number of protons, neutral or charged decay) and other interesting facts (ionization of tracks, etc.).

The records are fixed, unblocked and 40 character-long. This length for many reactions is over estimated.

The file is initiated by a deck of punched cards (scan-cards) carrying a few informations. The record is filled as soon as other informations are available in different steps of the chain.

When the direct access file is loaded, it is possible to obtain lists of all events concerning a reaction (lists of events to be measured). The program checking the goodness of measurements and preparing the input for reconstruction programs, updates at the same time the direct access file. At this point it is possible to obtain the first list of badly measured events; rates of rejected events are also available.

The good events comes to the reconstruction programs, and other informations update the direct access file. Misidentified events can be known in this step.

Different checks in various steps guarantee that no one event is processed more than once.

The successfull events are ready for the classification, that is done using the file's information and the data from reconstruction. The different information about same event are linked, and when they are compatible the event is identified. Doubtfull and misidentified events produce lists, that can be analyzed "by hand".

At the end of elaboration:

a) there is a DST updated with physical informations of identified events;

b) the direct access file is updated with informations concerning the classification of each event. This fact enables us to calculate automatically cross sections^(x).

In the case of automatic measuring devices, the direct access file can lead directly to the measuring device for chat concerns badly measured events.

(x) - In the initial loading of direct access file, a few record are reserved to record the number of tracks/photo in different photos.

5. - ESTIMATIONS FOR AN EXPERIMENT WITH A LARGE NUMBER OF EVENTS. -

The running system manages actually a pretty low number of events (20-30,000) with respect to the available in future experiments. It is possible to evaluate what will be necessary in the near future.

Since 100,000 events (2-6 prongs configurations), handled with 30 character/event blocked records in the direct access file, occupy nearly 10% of a 2314 IBM Disk in a IBM 360 system; then a large number of events ($\sim 10^6$) should occupy a 2314 Disk or a support of equivalent capacity. The cost of this equipment is very small in the overall cost of an experiment.

The time of elaboration for 100,000 events, requested in the interaction with the direct access file, is lower than 1% of total elaboration time.

The interactive subprograms are a few FORTRAN-routines inserted in usual management's programs (Thrifty⁽¹⁾/Thresh⁽²⁾, Grind⁽²⁾, Slake⁽¹⁾). Each of these routines occupies less than 5-6 K bytes of CPU (360/65 system). Fig. 3 displays the flux diagrams of insertion of different subroutines in the main chain.

6. - HISTOGRAMMING. -

Much time in physical data processing is requested for the histogramming of interesting physical data from DST. Generally a DST holds identified events of the same topology corresponding to one or more different reactions (e. g. $p3\pi - P3\pi + \pi^0$). Now two facts cause a considerable increase in computer time: a) the elaboration time in the step is directly related to the number of events on the DST; b) the same program needs a new reading of the DST, when the number of requested histograms exceeds a maximum number established by available computer storage.

The storage available is soon saturated when scatter plots are requested, especially when format recording is a matrix for each scatter plot. In this case a matrix of $n \times m$ order is needed, where n and m are the maximum number of channels for each dimension. Generally suitable codifying methods allows to allocate more than one channel in a word, but these methods put an upper limit to the number of events, which can be contained in a cell. This fact causes trouble for management of a large number of events.

Direct access devices allows the use of programs, which do not require a new reading of the tape, when the number of requested histograms exceeds the available core. A program of this kind is running in Bari⁽³⁾ on IBM 360/65 system.

PARTIAL DIAGRAM OF "THRIFTY" (PRE-THRESH PROGRAM)

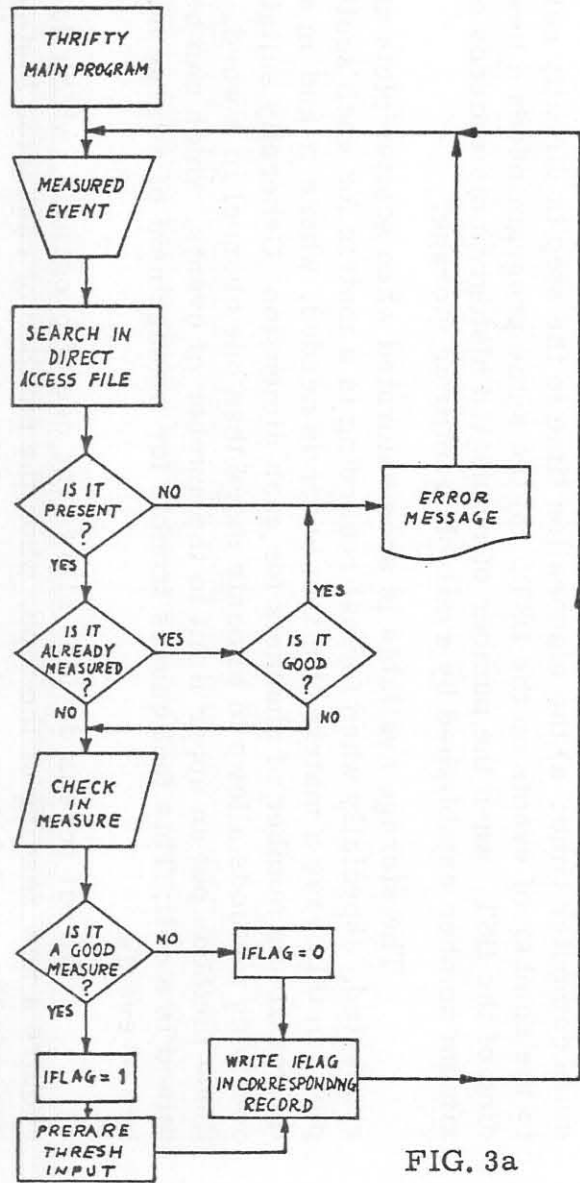


FIG. 3a

GRIND MODIFIED VERSION

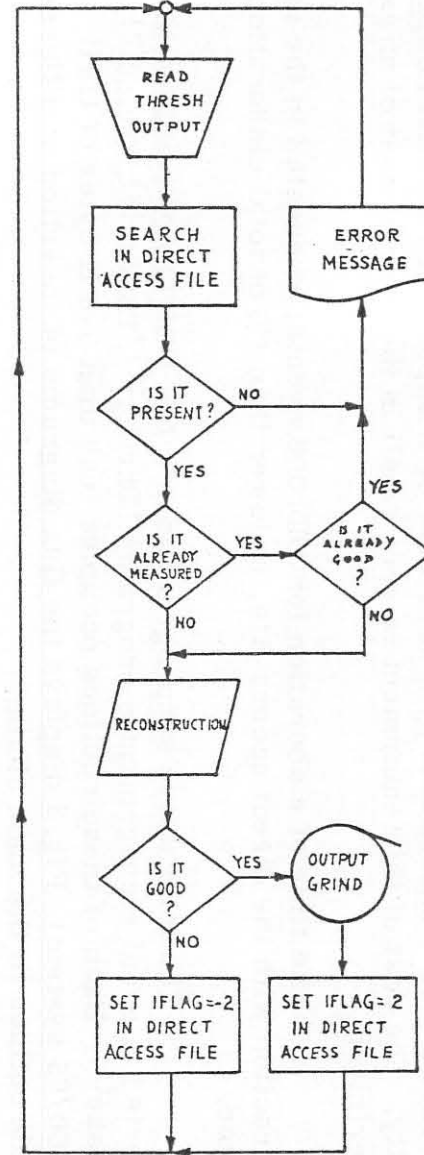
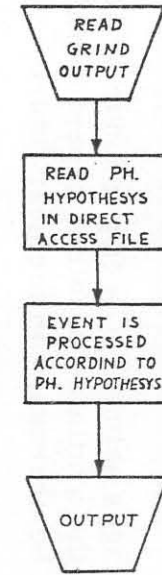


FIG. 3b

SLAKE MODIFIED VERSION



A large amount of information is recorded, transferring them in a direct access device (magnetic disk), when central memory is saturated. In order to effect this dynamical link, a list of single informations, fully in dependent, is built up. In our case this list is a set of n vectors of N dimensions for a monodimensional plot, or n' vectors of $2N$ dimensions for scatter plots; each component holds a single index in the first case or a pair of indexes in the second case. Indexes are the "coordinates" of the event in the plot. In such a way it is possible to obtain in a single run 100 monodimensional plots of 100 channels, and 50 scatter plots at 100×100 channels. For this run the "response time" (printing time is not included because it is independent on number of events) is due to the reading time of the DST; namely 5-6 minutes for 10.000 events, at 600 words/event (360/65 system). This method is acceptable for a few events, but becomes less and less usefull for large number of events, because the time required is just proportional to the number of events. Then a new philosophy is needed in this latter case for histogramming. This is described in the next paragraph.

7. - A NEW HISTOGRAMMING METHOD FOR LARGE NUMBERS OF EVENTS.-

The method is based on the hypothesis that the actual sequential structure of the DST is responsible of many troubles in handling large number of events.

The DST should be used only "una tantum" as general collection of information, while its present function could be replaced by a non-sequential structure of data; we call this structure DIF (Data Interactive File). Such a structure requests preferentially direct access supports, but also a configuration with a few sequential supports is tolerable, with a limitation in the efficiency of the method. A DIF consists of two files, we call them:

- a) Plot File;
- b) Test File.

The Plot file is organized in many subfiles, each one is directly addressable. Each subfile can be sequential. Each subfile holds n fixed blocks; each block contains m "cells", related to m events of DST; each celle holds k physical words of an event on DST. We will call this cell al so "multiplot". If N is the number of events of DST, then $n = N/m$.

The Test file consists of N logical records. Each logical record holds j words; all j words are to be considered as a logical unit, called "test-word". Then there are a multiplot and a test word for event. Also in the test file logical records (related with m events of the DST) are blocked in a physical record.

The k words of DST, related to physical values, belonging to a multiplot, must be selected taking into account the probability of correlations of physical words in a scatter-plot.

8.

Each multiplot has a potentiality of:

- k monodimensional histograms;
- $\binom{k}{2}$ bidimensional histograms;
- ⋮
- $\binom{k}{k}$ 1 k-dimensional histogram.

Each test word of the Test file has a length of $j \times N_{\text{bit}}$ bits, where N_{bit} is the number of bits in a word. The value 1 or 0 each bit means that a fixed condition in the event has TRUE or FALSE value.

The tests in the Test word must be defined "a priori" but these tests can be simple conditions; more elaborated tests can be assembled during the execution.

An optimal length for the test word is 5 word on a IBM/360 system (2, 5 on CDC 6600 system); this length corresponds to 160 simple tests, which if carefully selected allow repeated runs without other reading of DST.

8. - THE DIF BUILD-UP. -

The steps of the DIF's build-up are summarized below; (MULT(k,I), TEST(k) are transit memories for multiplots and tests respectively. Size is defined by:

I_{max} = total number of multiplot to be built

k_{max} = maximum central memory available. An optimal value for k_{max} is the track length of the direct-access device.

- 0 Set $k \leftarrow 1$; $ITIME \leftarrow 1$; MAX = maximum mass storage for each multiplot in direct access device;
- 1 Read one event from DST tape in TEMP; at end of file go to 6;
- 2 Transfer, for each I, M selected word from TEMP into MULT (k,I);
- 3 For each foreseen test assign true or false value and store it as a 1 or 0 bit in bit string TEST (k);
- 4 SET $k \leftarrow k + M$; is k equal to the maximum available storage? If yes go to 5, otherwise go to 1;
- 5 In a direct-access file write MULT bank at an address calculated for each multiplot as:
ADDRESS(I) = ITIME + MAX (I - 1);
on an other direct-access or sequential file (physically resident on an

other device) write TEST bank sequentially;

Set $ITIME \leftarrow ITIME + 1$; $K \leftarrow 1$; Go to 1;

6 Close the two files with a control-character;

Stop the run.

When the DIF is made, the histogramming starts. Generally a histogram is requested under a condition, which is structurally a logical "IF"; then the first step is to decodify this logical "IF" in elementary conditions, and to prepare a useful code for the next step.

The second step consists in the identification of subfiles, in which there is the histogramming word. In the case of a scatter plot there can be involved two words belonging to different multiplots, but this fact does not limit the efficiency of the method.

At this point a contemporary reading of the interesting multiplot and the test file permits the conditional rejection or acceptance of each event.

For what concerns:

a) the interpretation of logical "IF" in elementary condition; the method described in ref. (4) seems very suitable in bit string handling;

b) the identification of the subfile; this step can be solved by a cross-reference table, residing in a little file on the same direct access support. The cross reference table is transferred in the memory at each execution;

c) the reading of the interesting subfiles and test file is sequential. The histogramming method can be that just used in many programs.

A method to verify multiple conditions and cataloged procedures is described in ref. (4).

At the end the histogram can be displayed on printer or display by usual methods.

9. - TIMING FOR LARGE NUMBER OF EVENTS AND ON-LINE PROCESSING. -

An estimation of time to obtain one histogram by the method described in the previous paragraph has been attempted. In the case of more than one histogram the total time is nearly a multiple of estimated unitary time, because only few operations are independent from the number of histograms:

The estimation is for 100,000 events; due to the linear growth of the time with the number of events, the estimation for other number of events is immediate.

For 100.000 events each multiplot occupies 500.000 words, equivalent to near 15 cylinders on a 2314 dispak. The same occupation is needed for the test file at 160 test/testword. If the two data sets are on two different dispacks, the solar time for the elaboration is near 10 seconds. This time is essentially due to I/O; the elaboration time is much lower and it is largely contained in I/O time.

The traditional programs and a sequential DST involve the handling of many tapes for the same number of events. Then it is clear that an organization DIF type is very efficient, especially to obtain on line histogramming, requested by a terminal or by a video display.

The time is very small and compatible with a real time elaboration. For large number of events, too, the time in batch processing for an histogram alone is limited, then it is not necessary that all histogram are "a priori" forecasted.

It can be useful to point out that if on the plot file of DIF we record an index, or pair of indexes (see 1), instead of the words of DST, we can realize an occupation of memory, and a time of elaboration near 1/4 of that above estimated. This is due to the fact that the index values have a limited range (maximum value is the number of channels in a plot), and therefore many of them can be recorded in the same word. Naturally in this case the histogram steps are fixed.

This method allows to allocate 100.000 events in 75 tracks of a 2314 dispak, and the elaboration time for a histogram is 2-3 seconds.

REFERENCES. -

- (1) - M. Refice, THRIFTY USER'S MANUAL (Unpublished); B. Ghidini, SLAKE USER'S MANUAL (Unpublished).
- (2) - R. Bock, CERN Internal Report n. DD/EXD/62/10 (1962); CERN TC Program Library.
- (3) - N. Armenise, G. Piscitelli e A. Silvestri, PICAD: un programma di istogrammazione per sistemi con unità ad eccesso diretto, Rapporto INFN/TC-69/6 (1969).
- (4) - V. Capasso, A. Circella e A. Silvestri, Una logica a tre valori per il calcolo del valore di verità di funzioni booleane complesse, Rapporto interno.