

# THE NORTHEASTERN FOREST-INVENTORY DATA-PROCESSING SYSTEM I. INTRODUCTION



by  
**Robert W. Wilson Jr.**  
and **Robert C. Peters**

U. S. FOREST SERVICE RESEARCH PAPER NE-61  
1967

NORTHEASTERN FOREST EXPERIMENT STATION, UPPER DARBY, PA.  
FOREST SERVICE, U.S. DEPARTMENT OF AGRICULTURE  
RICHARD D. LANE, DIRECTOR

### **About the Authors . . .**

ROBERT W. WILSON JR. took his Bachelor's degree at The Pennsylvania State University in 1947 and his Master's at Yale University in 1948. He joined the U. S. Forest Service in 1948 and has worked in various research capacities for the Northeastern Forest Experiment Station. From 1961 to 1965 he was in charge of the Station's biometrics unit at New Haven, Conn. He is assigned at present to the Forest Insect and Disease Laboratory at West Haven, Conn.

ROBERT C. PETERS obtained his Bachelor's degree from the University of California in 1960 and his Master's at Yale University in 1961. He joined the Forest Service in 1961 as a research forester, and was assigned to the Station's biometrics unit from 1961 until 1965, when the unit was discontinued. Mr. Peters played a key role in the development of the data-processing system reported here.

# **THE NORTHEASTERN FOREST-INVENTORY DATA-PROCESSING SYSTEM**

## **I. INTRODUCTION**



### **CONTENTS**

A. Introduction .....	1
B. System outputs .....	4
C. Inventory data inputs .....	5
D. Compilation instruction inputs .....	5
E. System functions and organization .....	6
F. The EDIT subsystem .....	9
G. The SORT subsystem .....	9
H. The TABLE subsystem .....	11
I. The OUTPUT subsystem .....	14
J. System configurations .....	16

## PREFACE

THIS paper is the first in a series of ten papers prepared to describe the forest-inventory data-processing system of the Northeastern Forest Experiment Station. This system was devised for using modern, large-scale, high-speed computers in processing forest inventory data. The series will comprise the following papers:

- I. Introduction.
- II. Description of subsystem EDIT.
- III. Operation of subsystem EDIT.
- IV. Information for programmers—subsystem EDIT.
- V. Description of subsystem TABLE.
- VI. Operation of subsystem TABLE.
- VII. Information for programmers—subsystem TABLE.
- VIII. Description of subsystem OUTPUT.
- IX. Operation of subsystem OUTPUT.
- X. Information for programmers—subsystem OUTPUT.

## I-A. Introduction

ONE of the major projects of the U. S. Forest Service is a nationwide forest survey, which is designed to obtain useful and timely information about the timber resources of the United States. In the course of the surveys, which are made mainly on a state-by-state basis, great masses of detailed data are collected about timber volumes, growth, timber cut, and other characteristics of the timber resource.

In recent years the volume of information obtained from forest-survey field plots has increased greatly. The task of compiling and analyzing this mass of data with mechanical computing machines was both cumbersome and time-consuming.

A solution to this problem was seen in the development of the high-speed electronic computers. The Northeastern Forest Experiment Station, which was responsible for conducting the forest survey of the heavily forested Northeastern States, investigated the possibilities and devised the Northeastern Forest-Inventory Data-Processing System.

The purpose of this paper—and the nine companion papers that will follow—is to describe the system and its operation for the benefit of others who have similar problems in data processing.

Groundwork for the system was laid in 1960, when biometrician C. Allen Bickford<sup>1</sup> of the Northeastern Forest Experiment Station and biometrician W. G. O'Regan of the Pacific Southwest Forest and Range Experiment Station pointed out the possibilities and proposed the development of a standard computer program for the forest survey.

---

<sup>1</sup>Mr. Bickford, who retired from the U. S. Forest Service in 1963, is now on the faculty of the New York State University College of Forestry at Syracuse.

The authors of this paper were assigned to this project, and along with Bickford began to develop this program in the summer of 1961. The principal goals of the program were: (1) that the standard tables required of the forest survey were to be calculated and printed in a format suitable for publication; and (2) that each entry in every table would have a calculated sampling error.

In 1962, when this first program was developed, it was used for processing the inventory data collected in the forest survey of West Virginia. By this time several new inventories were under way, and some were completed except for data processing.

At this point it was evident that the standard computer program was already outdated because each new survey was different in at least one minor respect from the West Virginia survey; and although apparently minor, a difference often required a major alteration of the program. For this reason, the standard-program concept had to be changed to a more flexible one.

The present system is an outgrowth of several attempts to develop and to use both special-purpose and partially generalized systems to process the northeastern forest-survey inventory data. The features that give the system its flexibility are absolutely essential to an effective and economical data-processing operation. Differences between data-processing jobs appear to be inevitable. Even though all our processing jobs originate within a single organization operating under a stable set of general policies and objectives, differences of three principal kinds arise: (1) differences in required and optional resource statistics (output tables); (2) differences in measurement techniques, input data, and associated calculations; and (3) differences in sampling designs.

Differences of output tables arise primarily as reflections of differences of objectives among important users of resource statistics. In general, differences of this kind are small, but they are also unpredictable and difficult to adjust because they affect every step of the processing from input to output. Therefore provision must be made to accommodate these differences automatically from detailed descriptions of the output tables required for each processing job. Otherwise adjustments will be slow and expensive to make.

Differences in measurement techniques have many sources. Among the important ones are: technological improvements in measurement techniques, the value of data obtained with older techniques, and differences in the kind of resource being measured. Although organizational policy may properly demand standards of measurement, it would be both undesirable and unrealistic to expect that all differences can be eliminated. In fact, considering the number of reasonable alternatives to the measurement of each of the attributes of a resource, one must expect a large number of differences. The nature of these differences is largely unpredictable except in general terms. Virtually the only way that these differences can be accommodated is by providing for the input of detailed descriptions of measurement checks and calculations with the data for each processing job.

Although differences in sampling designs may be large in effect, the number of different designs in common use is few; and the associated differences in compilation procedures are relatively small. Consequently it has been possible to incorporate in the system fully detailed procedures for each of the more common alternatives. Only the appropriate alternative need be indicated in the specifications for each processing job.

The system as finally developed encompasses all the processing of inventory data from field recording sheets to the printing of labeled tables of resource statistics. The system features variable inputs, variable outputs, and variable processing procedures for getting from one to the other. The user needs only to specify his particular processing problem in terms of the procedural options offered in the system. Consequently, the system is a general one, applicable to a wide variety of inventory-processing problems. Yet, for the user who has adopted standard procedures that he uses repeatedly, the system retains the advantage of a special-purpose system. Once he has exercised his options, he can incorporate them into his system, and they remain a part of it until such time as he wishes to change it. Changes are accomplished simply by making appropriate changes in the choice of options; there is no need to develop a completely new system or to undertake the hazardous task of patching the original system.

In short, the system provides all the information channels and detailed procedures necessary for processing inventory data while giving the user freedom to choose how his data should be handled for best results.

Other features of the system are its organization, designed to require the minimum of manual data handling consistent with the desired flexibility; and its use of large modern high-speed computers. The result is a capacity for consistent, rapid, low-cost processing of large volumes of inventory data.

### **I-B. System Outputs**

The primary output from the system consists of a set of up to 40 one- or two-dimensional tables, designed to contain the required resource statistics for each survey unit or group of survey units for which inventory data inputs are given. Similar outputs may be obtained for subdivisions of survey units. If more than 40 output tables are required, the data must be passed through at least part of the system a second time; otherwise only one pass of the data is required.

Each table contains a particular summary of one or more sampling-unit attributes over the entire population of sampling units. The basic statistic may be either the sums or the means of the sampling-unit attributes. If the inventory data input is a sample from the population, additional tables of the same format are output automatically; they contain the variance of the statistic, and optionally, its standard error. Row, column, and grand totals are also provided automatically for each table. The tables are printed complete with appropriate labels provided by the user.

In addition to the primary output, almost any other type of output can be obtained, given adequate specifications by the user. The intermediate outputs from the main processing stages can be obtained simply by specifying that they be saved. Special outputs representing each unit record and/or any sets of unit records can be obtained from one stage of processing, EDIT. Special outputs representing each sampling unit can be obtained from another processing stage, TABLE.

### **I-C. Inventory Data Inputs**

Inventory data are nothing more than a catalog or itemized list of objects, or categories of objects, that have actually been observed. It is in the form of such a list that data are used in processing, but the inputs need not be presented to the system in this form. The only limitations on the data inputs are that each item in such a list is expressible as a unit record that contains complete identification and description (values for each relevant attribute) of the item in 132 or fewer data fields, and that the formats of each unit record be identical.

In forest inventories the relation between the objects actually observed and the statistical sampling units may vary from one sampling plan to another. Generally, the statistical sampling unit is a plot of ground of given dimensions, whereas at least some of the sampling-unit attributes can be actually observed only on a subdivision of the sampling unit such as a point or a tree. Consequently, the input data for a sampling unit (the object of observation for purposes of compilation) are a set of unit records in which the actual observations on each subdivision of the sampling unit are recorded. The set may contain any number of unit records not exceeding 160. Other relations between objects actually observed and statistical sampling units reduce to special cases of the general case just described and can be handled accordingly. In short, the inventory data input is always handled as a simple itemized list of observed objects. What these objects are and how the data will be processed depend upon the sampling plan and field procedures under which the data were obtained.

### **I-D. Compilation Instruction Inputs**

The compilation instructions describe the particular data-processing job at hand. The instruction inputs are treated by the system as data, but they are a special kind of data. They are incorporated into the system in such a way as to complete and/or modify it as required for the job to be done. To a large extent, this operation is done automatically. Because of the way they are handled, the instruction inputs may be considered as an integral

part of the system once they have been presented to it, but they may also be changed at will if the need arises.

The instruction inputs describe: (1) the outputs required of the system; (2) the data inputs to the system; (3) the measurement checks and calculations required to obtain values of all relevant attributes of the sampling units; (4) the rules whereby the set of unit records representing the sampling unit are transformed to a set of tables for the sampling unit that have the same format as the output tables; and (5) the sampling design and estimating procedures to be followed in processing.

### **I-E. System Functions and Organization**

The system is designed to produce output tables of resource statistics from the following types of inventory:

1. Complete (100-percent) coverage of sampling units in a survey unit.
2. Simple random sample of sampling units in a survey unit.
3. Systematic sample of sampling units in a survey unit, with random elements in sample selection.
4. Stratified random sample of sampling units in a survey unit, with known stratum weights.
5. Stratified random sample of sampling units in a survey unit, with stratum weights estimated from a primary sample.

Direct estimates of the resource statistics can be made from each of these inventory types. For the fifth type of inventory, several alternative means of indirect (ratio and regression) estimation are also provided. These have been designed primarily to accommodate the northeastern forest survey and will be described more fully in the discussion of that survey in part X.

The operations needed to make these estimates are grouped into subsystems, each covering a major phase of the processing. The result is an orderly sequential flow of information from the user to the system, through the successive processing steps and back to the user, as illustrated by the case of processing a single sample (fig. 1.) This configuration of the system is the normal one for a processing job, but variations of this arrangement may also be

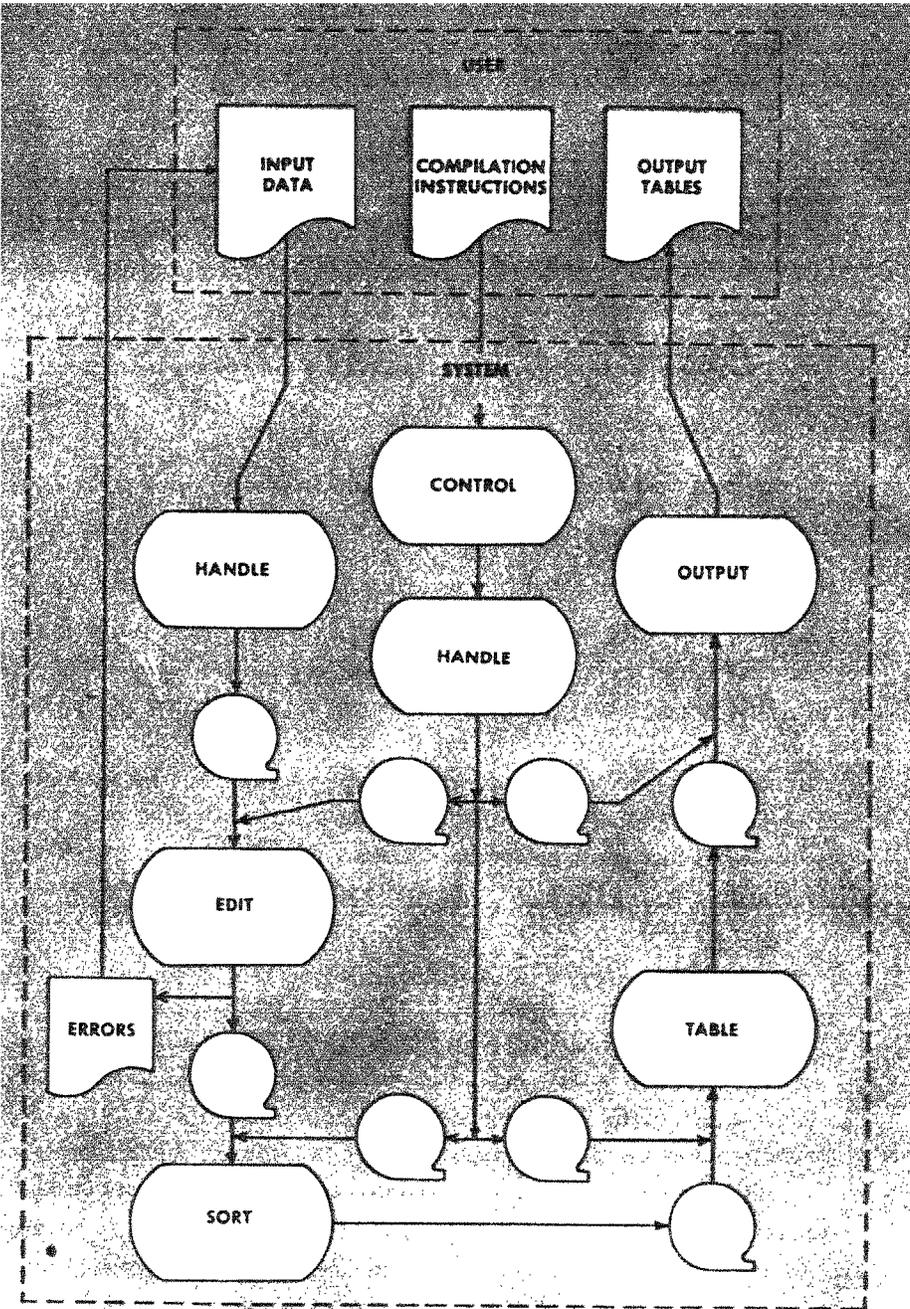


Figure 1.—System configuration for processing a single sample.

employed. System configurations will be discussed directly after the discussion of the functions of each subsystem (I-J).

The functions of the CONTROL subsystem are to accept compilation instructions from the user, to translate the instructions into appropriate system configurations and control data for each stage of processing, and to monitor and supervise the entire processing job. The only component of this subsystem is a person who is thoroughly acquainted with forest-inventory procedures and with the operations of this system.

The principal function of the HANDLE subsystem is to handle information (both inventory data and control data) within the system. It performs such operations as the punching, verifying, rearranging, and taping of information for later use, and the storage and retrieval of such information in master files. The components of the subsystem are: (1) people who operate machines and perform clerical work; (2) hardware such as key punches, verifiers, and an IBM 1401 System; and (3) software (computer programs) for the latter, to be used in maintaining the master files.

It is important to note that these two subsystems are presented as conceptual rather than as matters of fact. The functions encompassed by these subsystems must be performed, but because the principal components are people whose salaries must be paid and machines that must be rented, the user or his organization is responsible to set them up. The particular configurations of these systems will depend upon the volume of work, the cost limitations, and so forth.

The remaining subsystems are composed primarily of computer programs. The heart of each subsystem is a computer program tailored to the needs of the system and bearing the same name as the subsystem. The programs have been written in FORTRAN IV for use with Yale University's IBM 7094/7040 Direct Coupled Computer System.<sup>2</sup> The subsystems are operational under the IBSYS Monitor and will run with little or no modification on other computers that accept FORTRAN IV, have binary arithme-

---

<sup>2</sup> Mention of a particular product should not be construed as an endorsement by the Forest Service or the U. S. Department of Agriculture.

tic, a 32K core (minimum of 36 bits per word), and 5-tape drives or equivalent input-output devices.

### **I-F. The EDIT Subsystem**

The EDIT subsystem functions as a unit-record processor (fig. 2). It processes one record at a time from the data input list previously described. Its functions are to check the values in data fields or groups of data fields in each record for validity or reasonableness, to generate values for additional data fields from input values of other data fields, and to provide summaries of sets of records consisting of the sums over the sets of values in one or more data fields.

These functions are performed by a set of pre-programmed operations. For each check or calculation to be performed on a record, the user selects an appropriate operation, designates the data fields to which it is to be applied, and furnishes the numerical operators required.

The outputs from the subsystem are of three kinds: (1) the correct, augmented records, which are put on a magnetic tape file; (2) error messages identifying incorrect records and on which the types of error are printed; and (3) data-field summaries, which are both printed and punched in the form of tables. When error records have been corrected, they are passed through the subsystem a second time and, if they pass the checks, they are added to the correct record tape file and the data field summary outputs.

### **I-G. The SORT Subsystem**

The SORT subsystem is a general IBM system which is used to arrange the correct record output from the EDIT subsystem in the order required for subsequent processing. The records for a sampling unit are arranged in sets and ordered within the sets, if necessary. These sets are in turn arranged into groups for each sample stratum and ordered within the groups if necessary. The sample stratum groups are, in turn, grouped into survey unit groups (arbitrary subdivisions of a geographical area) which, in turn, are ordered within the geographical area. This exhausts the hierar-

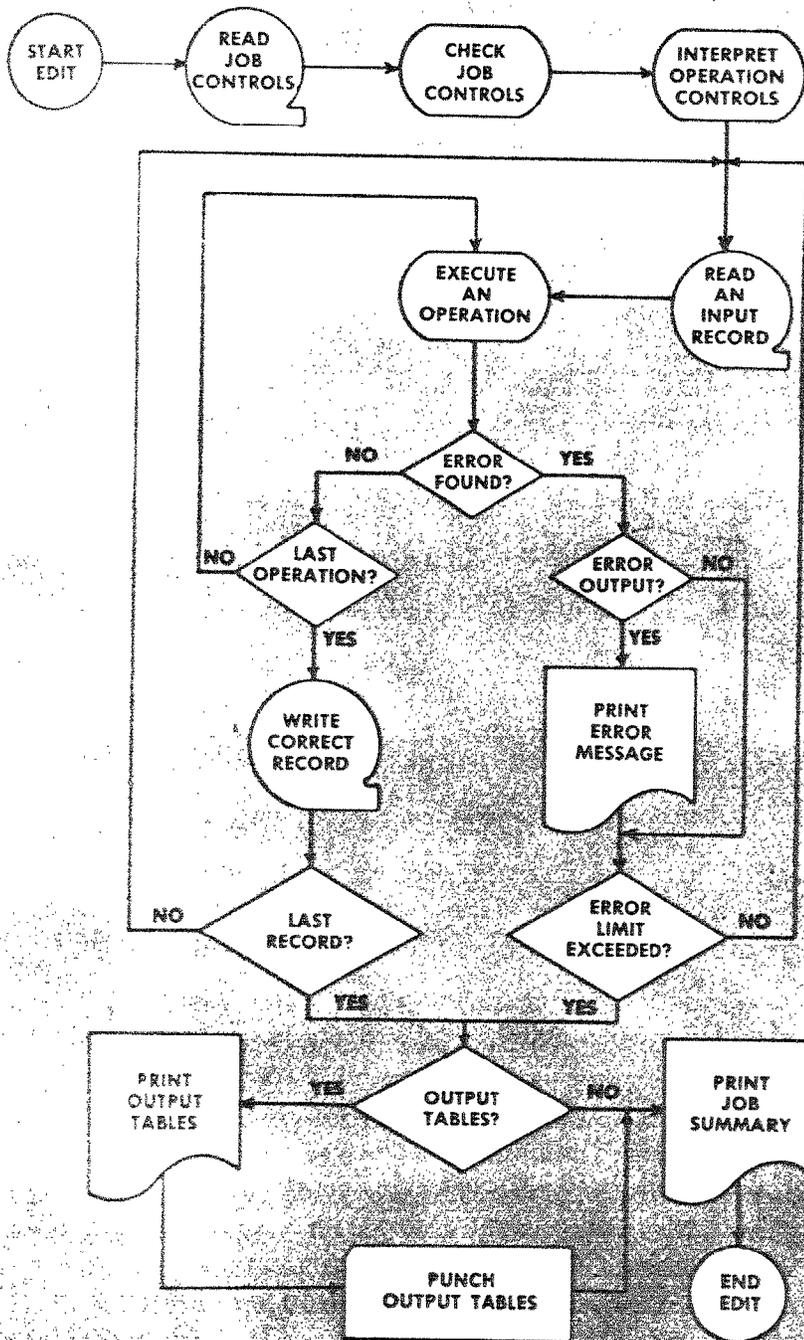


Figure 2.—A generalized flow chart of EDIT.

chy of records that are available in the system. The levels that are recognized depend upon the type of inventory.

### **I-H. The TABLE Subsystem**

The TABLE subsystem has three primary functions: (1) to transform the data input for each sampling unit into a set of tables with the same format as the final output tables; (2) to sum these sets of tables, cell by cell, over all the sampling units in the sampling-unit set; and (3) to compute output statistics for the tables of each sampling-unit set (fig. 3). Four options, or alternative combinations of output statistics, are available:

1. Sums of tabular sampling-unit values over all sampling units in the sampling-unit set.
2. Arithmetic means of tabular sampling-unit values over all sampling units in the set of sampling units.
3. Arithmetic means and their variances of tabular sampling-unit values over all sampling units in the sampling-unit set.
4. Arithmetic means, their variances, and covariances between the means of each table cell and the mean over all cells in the tables.

The last option is designed primarily for the northeastern forest survey to use in estimating the variance of ratio estimates. (See part X in this series).

The data inputs to this subsystem are simply the ordered sets of sampling-unit inputs, each consisting of an ordered set of unit records, as described under the SORT subsystem.

Processing proceeds one sampling unit at a time. The critical function of transforming the sampling-unit input data to output tables is performed by a set of pre-programmed operations. The required tables are formed sequentially for each sampling unit. In general, table entries are made in the following manner: The value in some data field (designated in the compilation instructions) is entered into the table at a location determined as a function of the values in one or two other data fields of that record. The functions and the data fields to be used in determining the location of an entry are also obtained from the compilation instructions.

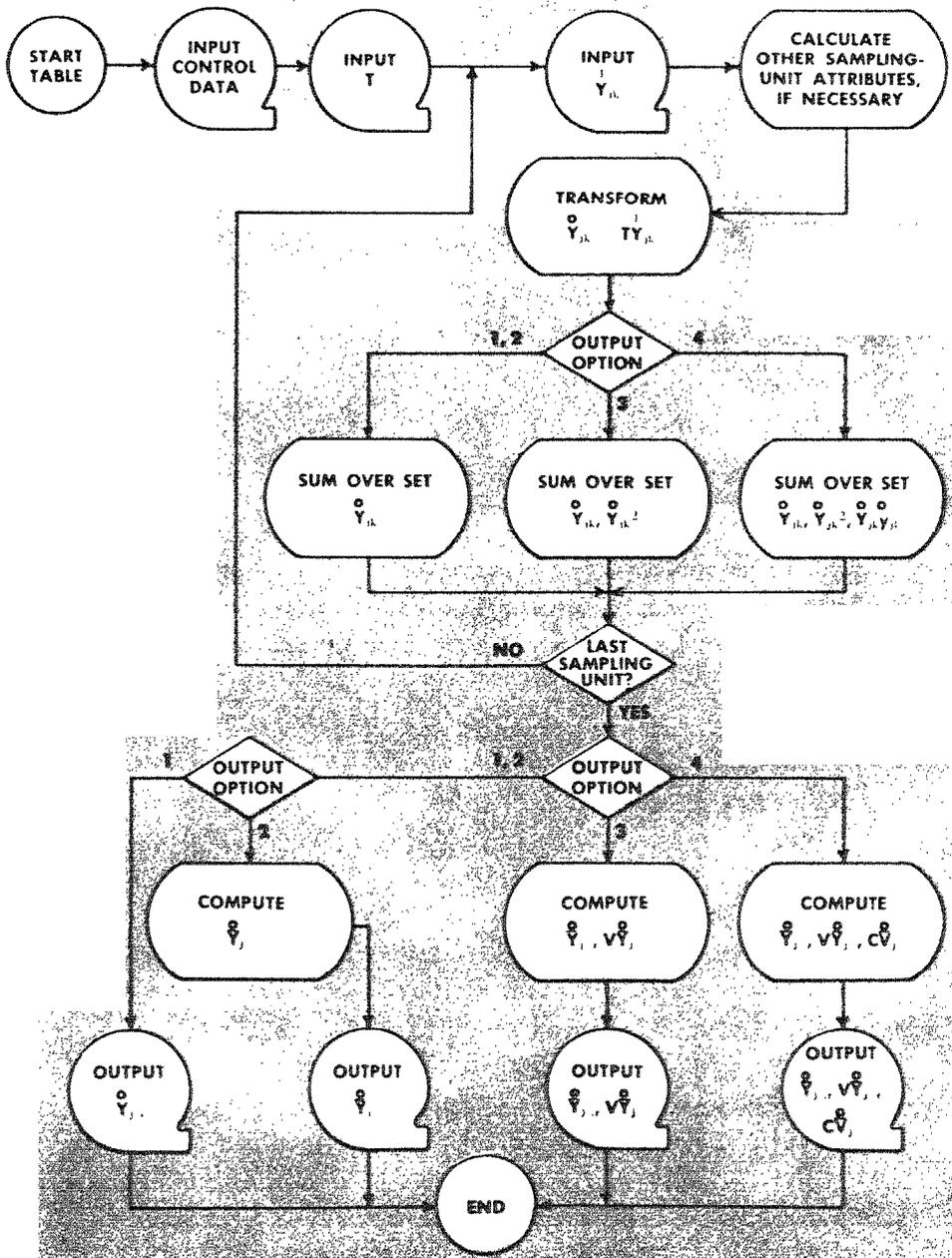


Figure 3.—A generalized flow chart of TABLE.

$i$  = Subscript for the  $i$ th element of a vector.

$j$  = Subscript for the  $j$ th sample stratum.

$k$  = Subscript for the  $k$ th sampling unit.

$T$  = A set of rules (transform) whereby the elements of the sampling-unit attribute vectors (inputs),  $Y_{jk}^I$ , are redistributed to form the sampling-unit attribute vectors (output tables),  $Y_{jk}^O$ .

$Y_{jk}^I$  = An input vector containing a sampling-unit attribute (a data field from the input data matrix for a sampling unit).

$Y_{jk}^O$  = A sampling-unit attribute vector (output table) which represents a summary of the input vector,  $Y_{jk}^I$ , for that attribute.

$y_{.jk}^O$  = The sum (total inventory) of the elements of the sampling-unit attribute vector (output table),  $Y_{jk}^O$ .

$Y_{j.}^O$  = An attribute output vector (table), containing the arithmetic means for the sampling-unit attribute vectors,  $Y_{jk}^O$ , for the stratum.

$VY_{j.}^O$  = The variance of  $Y_{j.}^O$ .

$CV_{j.}^O$  = An output vector (table), containing the mean covariances of  $Y_{jk}^O$  with  $y_{.jk}^O$ .

Provision is also made in the subsystem for excluding certain unit records from the tables, and for making more than one entry in a table for each unit record.

### **1-1. The OUTPUT Subsystem**

The primary function of the OUTPUT subsystem is to produce the complete, fully labeled output tables of resource statistics describing the population of sampling units in each survey unit given as input (fig. 4). There are six ways in which these tables may be produced:

1. Tables for one or more sets of sampling units that represent the survey unit may be read in, added together, and printed. The input for this alternative is the output from the first option of the TABLE subsystem (or its equivalent).
2. Tables for a set of sampling units representing the survey unit may be read in, multiplied by an appropriate expansion factor, and printed, with standard errors. The standard errors are computed for a simple random sample. The input for this alternative is the output from the third option and the expansion factor.
3. Tables for each sample stratum within the survey unit may be read in, weighted by stratum size, summed, multiplied by appropriate expansion factors, and printed, with standard errors. The sampling errors are computed for a stratified random sample with known stratum weights. The inputs for this alternative are the output from the third option, the stratum weights, and the expansion factors.
4. Tables for each sample stratum within the survey unit may be read in, weighted by stratum size, summed, multiplied by appropriate expansion factors, and printed, with standard errors. The standard errors are computed for a stratified random sample with stratum weights estimated from primary sample. The inputs for this alternative are the outputs from the third option, the estimated stratum weights, and the expansion factors.
5. Tables for each sample stratum within the survey unit are read in, converted to ratios (ratios of the values in the total cells of

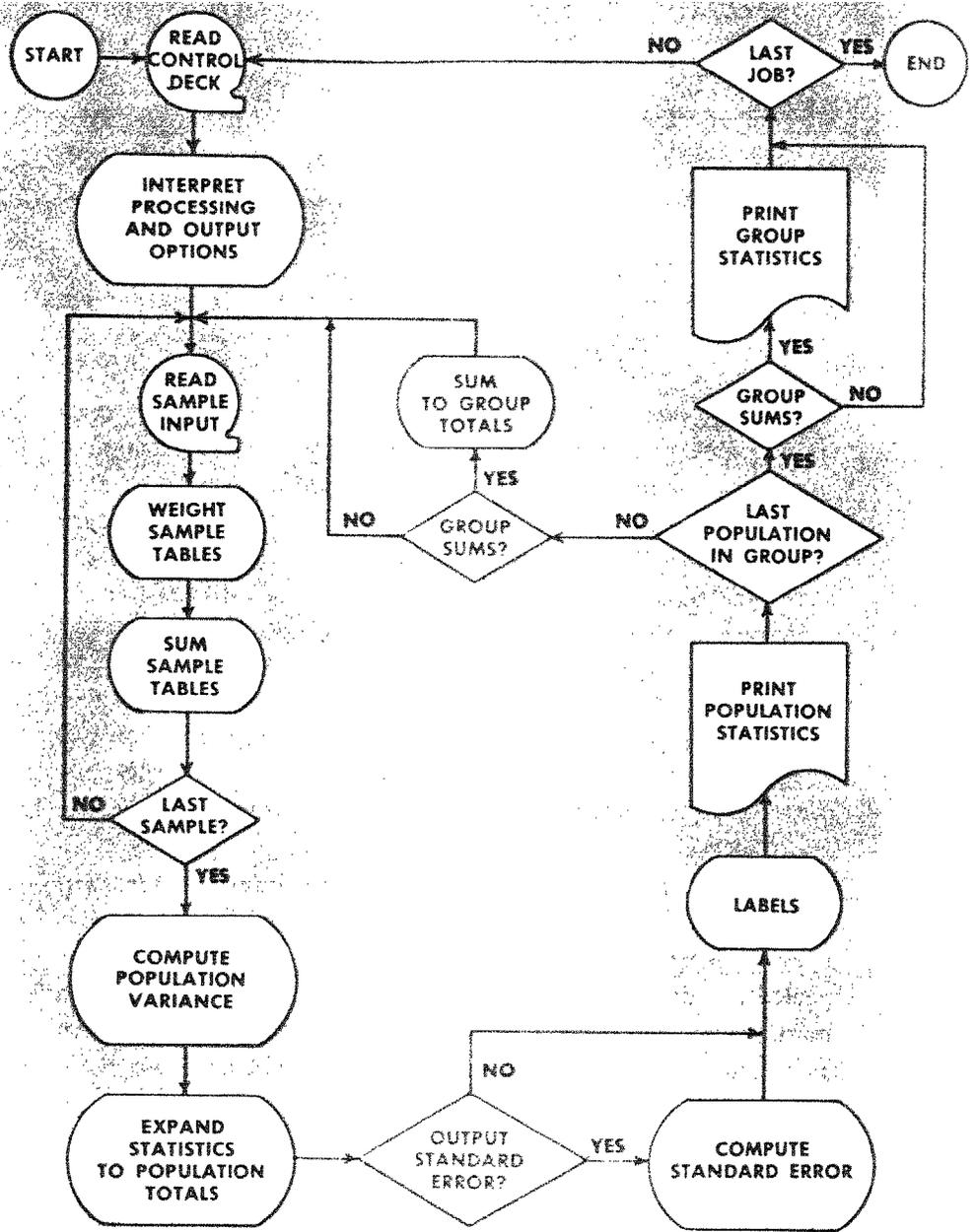


Figure 4.—A generalized flow chart of OUTPUT.

the tables), multiplied by independent estimates of the table totals, weighted by stratum size, summed, multiplied by appropriate expansion factors, and printed, along with standard errors. The standard errors are computed for a separate stratified ratio estimate, including estimated stratum weights. The inputs for this alternative are the outputs from the fourth option, the stratum weights, the stratum totals for each table, and appropriate expansion factors.

6. Tables for each sample stratum within the survey unit are read in, weighted by stratum size, summed over strata, converted to ratios (ratios of the values in the total cells of the tables), multiplied by independent estimates of table totals and appropriate expansion factors, and printed, with standard errors. The standard errors are computed for a combined stratified ratio estimate, with estimated stratum weights. The inputs for this alternative are the outputs from the fourth option, the estimated stratum weights, the independent estimates of the survey-unit totals, and the expansion factors.

In addition to the survey-unit outputs just described, two other kinds of outputs can be obtained from this subsystem. When two or more survey units are processed successively in the same pass, the table estimates and their variances can be accumulated over the group of survey units, standard errors for the group can be computed, and the results for the group can be printed in the same format as the survey-unit results. Alternatively, when only one survey unit is processed in a pass, the table inputs can be used repeatedly with different weights, expansion factors, etc., to obtain estimates and sampling errors for any number of subdivisions such as counties of the survey unit. Again, the results are printed in the same format as the survey-unit results.

## **I-J. System Configurations**

The system as described so far will accommodate most forest-inventory data-processing jobs. In some cases it may be possible to

bypass some subsystems, such as EDIT and/or SORT. In some cases the data may be partially processed by other means, and only one subsystem may be needed. Either of these situations can be met, provided the necessary inputs are available. Any of the processing subsystems except OUTPUT can be used independently of the system as a whole.

On the other hand, some processing jobs may require more elaborate configurations of the system, usually based on multiples of the configuration shown in figure 1. For example, the northeastern forest survey uses three stratified ground samples in making its estimates, all with stratum weights estimated from primary photo samples.

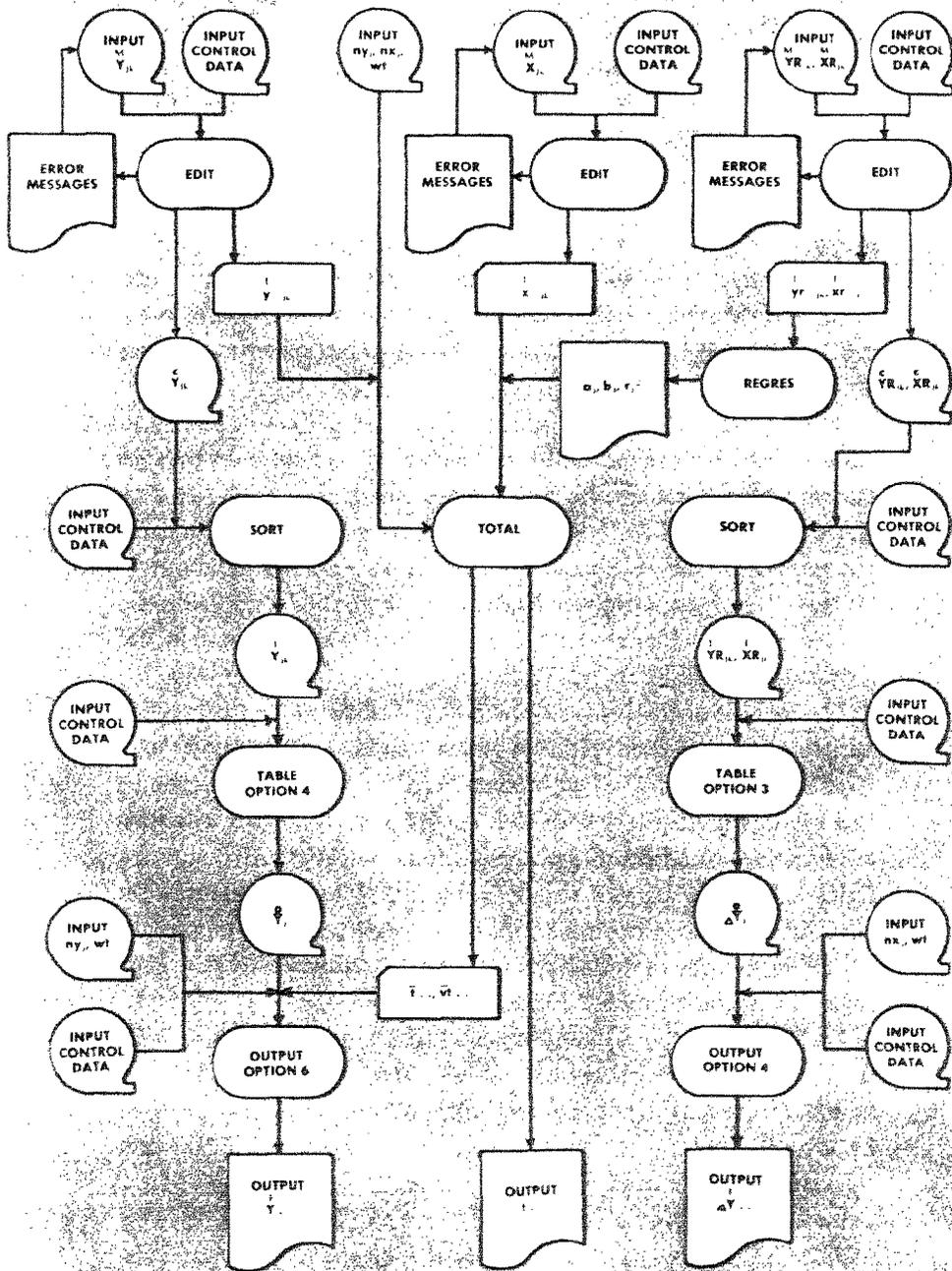
A current sample is used to estimate total areas and total volumes, and to estimate ratios for the distribution of total inventories into the output tables. As already described, the ratios are estimated in the OUTPUT subsystem by either of two procedures: the ratios may be estimated stratum by stratum (separate ratio estimates) or from weighted means over all strata (combined ratio estimates). The latter procedure is currently being used in the northeastern forest survey.

The sample of the previous inventory is updated by a regression sample of remeasured plots to obtain another estimate of total areas and volumes. The two estimates of totals are then combined and used as multipliers with the tabular ratio estimates to obtain the final output tables of area and volume statistics. Although the ratios and totals are not strictly independent in this example the effect of covariance on the sampling errors is assumed negligible.

Finally, output tables of growth and mortality statistics are estimated directly from the regression sample.

Consequently, the system configuration for processing the northeastern forest survey (fig. 5) provides that each of two samples be passed through the processing steps shown in figure 1. It will be noted also that there are two additional subsystems in this configuration that are used to produce the estimates of total areas and volumes.

The TOTAL subsystem provides the means to produce independent estimates of the grand totals of each output table (max-



**Figure 5.—System configuration for Northeastern Forest Survey processing.**

- $i$  = Subscript for the  $i$ th element of a vector.
- $j$  = Subscript for the  $j$ th sample stratum.
- $k$  = Subscript for the  $k$ th sampling unit.
- $Y_{jk}^M$  = Mixed (unsorted) sets of sampling-unit data input records from the current sample.
- $Y_{jk}^C$  = Correct but unsorted sets of sampling-unit data input records from the current sample.
- $Y_{jk}^I$  = Correct and sorted sets of sampling-unit data input records from the current sample.
- $Y_{j.}^O$  = Sets of output tables containing means and their variances over sets of sampling units of sampling-unit attributes from the current sample.
- $Y_{j.}^F$  = Sets of output tables containing the area and volume statistics required, with sampling errors.
- $t_{j.}$  = A set of final estimates of output table totals, with sampling errors.
- $\Delta Y_{j.}^F$  = Sets of output tables containing the growth and mortality statistics required, with sampling errors.
- $X_{jk}$  = Sets of sampling-unit data input records from the initial inventory sample.
- $YR_{jk}, XR_{jk}$  = Paired sets of sampling-unit data input records from the regression sample (re-measured plots).
- $Y_{j.k}^I$  = Set of output-table totals for each sampling unit.
- $x_{j.k}$  = Set of values of independent regression variables corresponding to output-table totals for each sampling unit.
- $ny_j, nx_j$  = Sets of sample stratum weights from the primary samples.
- $wt$  = Expansion factor or size of survey unit (total number of sampling units).

imum of 40 tables) from stratified random samples, with stratum weights estimated from primary samples. Although these are intended primarily for use with the tabular ratio estimates, they may also be useful in their own right when full output tables are not required or when estimates of total inventory are needed very quickly. Four alternatives are available, each with a corresponding set of sampling errors:

1. Separate stratified regression estimates.
2. Combined stratified regression estimates.
3. Stratified direct estimates.
4. Combination of two independent estimates made by any of the previous alternatives.

The data inputs are the data-field summaries of volumes and area classes for each sampling unit of a sample obtained as special output from the EDIT pass of each sample. Either one or two samples may be processed at one time, but the resulting estimates can be combined only if the two independent estimates are processed in one pass.

The REGRES<sup>3</sup> subsystem provides the means to obtain stratum-by-stratum estimates of regression parameters and/or ratios, and their sampling errors for the totals of up to 40 output tables. The data inputs are the data-field summaries of paired volumes and area classes for each sampling unit in the regression sample (re-measured plots) obtained as special output from the EDIT subsystem.

---

<sup>3</sup> REGRES: an abbreviation for "regression sample."

