

Historic, archived document

Do not assume content reflects current scientific knowledge, policies, or practices.

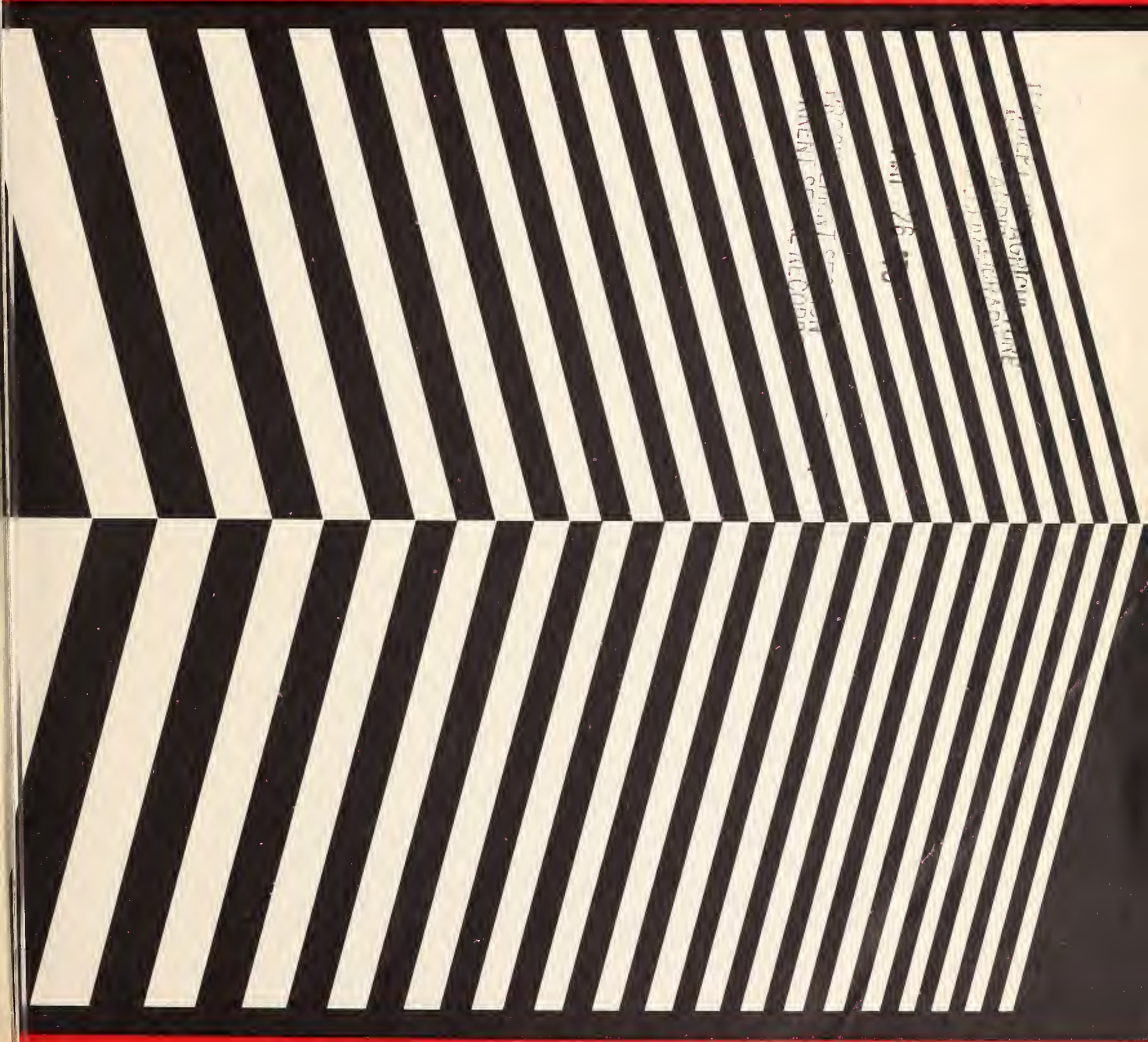
serve.
9
892Esc

ANALYZING IMPACTS OF EXTENSION PROGRAMS

EXTENSION SERVICE

U.S. DEPARTMENT
OF AGRICULTURE

ESC-575



LIBRARY
U.S. DEPARTMENT OF AGRICULTURE
WASHINGTON, D.C. 20250
APR 26 1968

Extension's productivity and accountability can be greatly advanced through evaluation of its programs. This publication is intended to help Extension administrators, program leaders, and specialists to fulfill their program evaluation responsibilities.

Emphasis is placed on maximizing the usefulness of program evaluations in decisionmaking on program priorities and modifications.

Harmonizing ideal program evaluation with available resources is a challenge; program evaluation, especially the collection of evidence, can be costly. *Analyzing Impacts of Extension Programs* presents optional levels of evidence with varying degrees of expense.

Some of the ideas in this publication are yet to be verified, but are presented to guide thinking about Extension program evaluation and to lead to tested principles of evaluation.

Administrative support for preparation of this publication was extended by C. A. Williams, Deputy Administrator, Program and Staff Development, Extension Service, U.S. Department of Agriculture.

The Extension Service of the U.S. Department of Agriculture offers its programs to all eligible persons regardless of race, color or national origin, and is an Equal Opportunity Employer.

Cooperative Extension Work: U.S. Department of Agriculture and State Land-Grant Universities Cooperating. Issued April 1976.

CONTENTS

Abstract	2
Introduction	3
A Chain of Events in Extension Programs	3
A Hierarchy for Program Evaluation	5
Reaching Program Objectives	7
Selection of Level of Evidence	8
A Pyramid of Evidence for Program Evaluation	9
Evaluation Criteria and Quality of Evidence	11
Proxy Measures	15
Designs for Identifying Source of Impact	15
The Field Experiment	16
Matched Set Design	16
Time-Trend Studies	18
“Before-After” Study	18
The Survey	19
The Case Study	20
Using and Appraising Evaluation Studies	20
Summary and Conclusions	21
Acknowledgments	21

This publication presents a framework, guidelines, strategy, and methods for evaluating Extension education programs. Extension programs are viewed in terms of seven levels of objectives and evaluative evidence: (1) inputs, (2) activities, (3) people involvement, (4) reactions, (5) change of knowledge, attitudes, skills, and/or aspirations (KASA), (6) practice change, and (7) end results.

Levels 1 and 2 characterize Extension's efforts. Level 3 includes the people involved by Extension and the nature of their involvement; levels 4 through 7 cover the responses by these people and others. Responses range from the immediate and direct to the long-term and indirect consequences of Extension's actions.

The foregoing levels vary in: (a) the extent to which they can provide evidence of Extension's impact and (b) the amount of resources required for obtaining evidence. Evidence of Extension program impact becomes stronger in ascending the levels. However, obtaining evidence at higher levels generally requires more evaluative resources. The level(s) of evidence chosen for a particular program evaluation will vary with the decisions it is to assist, the nature of the program, and the circumstances of its evaluation. Proxy indicators are suggested, in order to maximize strength of evidence in lower cost assessments of Extension's effectiveness.

Program evaluations may be relied upon to assist decisionmaking to the extent that they provide high-quality evidence of accomplishment of program objectives and identify Extension's extent of contribution to such accomplishments.

ANALYZING IMPACTS OF EXTENSION PROGRAMS

INTRODUCTION

"Are Extension programs succeeding?" is a question asked frequently by officials at all levels of Government, legislators, university administrators, and Extension workers themselves. This publication provides guidance in evaluating Cooperative Extension education programs.¹

Judgments about program effectiveness will be made one way or another. However, there is mounting demand by legislators, policymakers, and administrators that program effectiveness be demonstrated through formal evaluations. These demands reinforce the desire by Extension staff to obtain sound evidence of the extent to which Extension programs are successful. *Formal evaluation entails conscious procedures for placing value on programs according to (1) explicit criteria and (2) designs for collection and analysis of evidence.*

Program evaluation is part of the overall program development process, which includes: (1) identifying problems and selecting long-range objectives; (2) specifying these objectives and the strategy, activities, and budget designed to achieve them; (3) conducting activities; (4) evaluating the program's strategy and impact; and (5) using this evaluation along with other information in subsequent program development.

Impact evaluation is assessment of a program's effectiveness in achieving its ultimate objectives or assessment of relative effectiveness of two or more programs in meeting common ultimate objectives.²

¹ Cooperative Extension education is defined herein as noncredit individual, group, and mass instruction directed toward practical problem-solving. Usually conducted off-campus and informally, Cooperative Extension programs are an outreach of Land-Grant Universities and Colleges. Cooperative Extension Service programs are generally mutually funded and directed by local, State and national sources. See U.S. Department of Agriculture, National Association of State Universities and Land Grant Colleges Study Committee, *A People and A Spirit*, Fort Collins, Colo., Colorado State University, 1968.

² Scriven, Michael, "The Methodology of Evaluation," *Perspectives of Curriculum Evaluation*, Ralph Tyler, Robert Gagne, and Michael Scriven (eds.), pp. 39-83, Chicago, Ill., Rand McNally, 1967.

Stufflebeam, Daniel L., "Toward a Science of Education

By Claude F. Bennett
Specialist, Educational Methodology and Evaluation
Program and Staff Development

The major purpose of program evaluation is to assist in reaching decisions on future directions, design, and funding of programs.³ Decisions on whether programs should be terminated, curtailed, maintained, or expanded are aided by program evaluations.

Such evaluations may also suggest reformulation of program objectives, strategy, delivery organization, educational methodology, and intended audiences.

This publication identifies seven broad categories of criteria which are useful in formally evaluating the effectiveness of Extension programs and attempts to provide guidance in choosing evidence regarding these categories.

A CHAIN OF EVENTS IN EXTENSION PROGRAMS

Figure 1 shows a "chain of events" assumed to characterize most programs of Extension education. Although the events selected oversimplify reality, they provide a "mind-hold" on Extension programs. The events chart the behavior of both Extension and the people involved in its programs.⁴

"Inputs" (lower left of fig. 1) are selected on the assumption that problem solution may require resource expenditures. With these inputs, "Activities" can be performed; e.g., publicizing programs or "putting across" educational content.

Activities "Involve People" (participants) who have "Reactions", i.e., some degree of interest

Evaluation," *Educational Technology* 8 (July 1968), pp. 5-12.

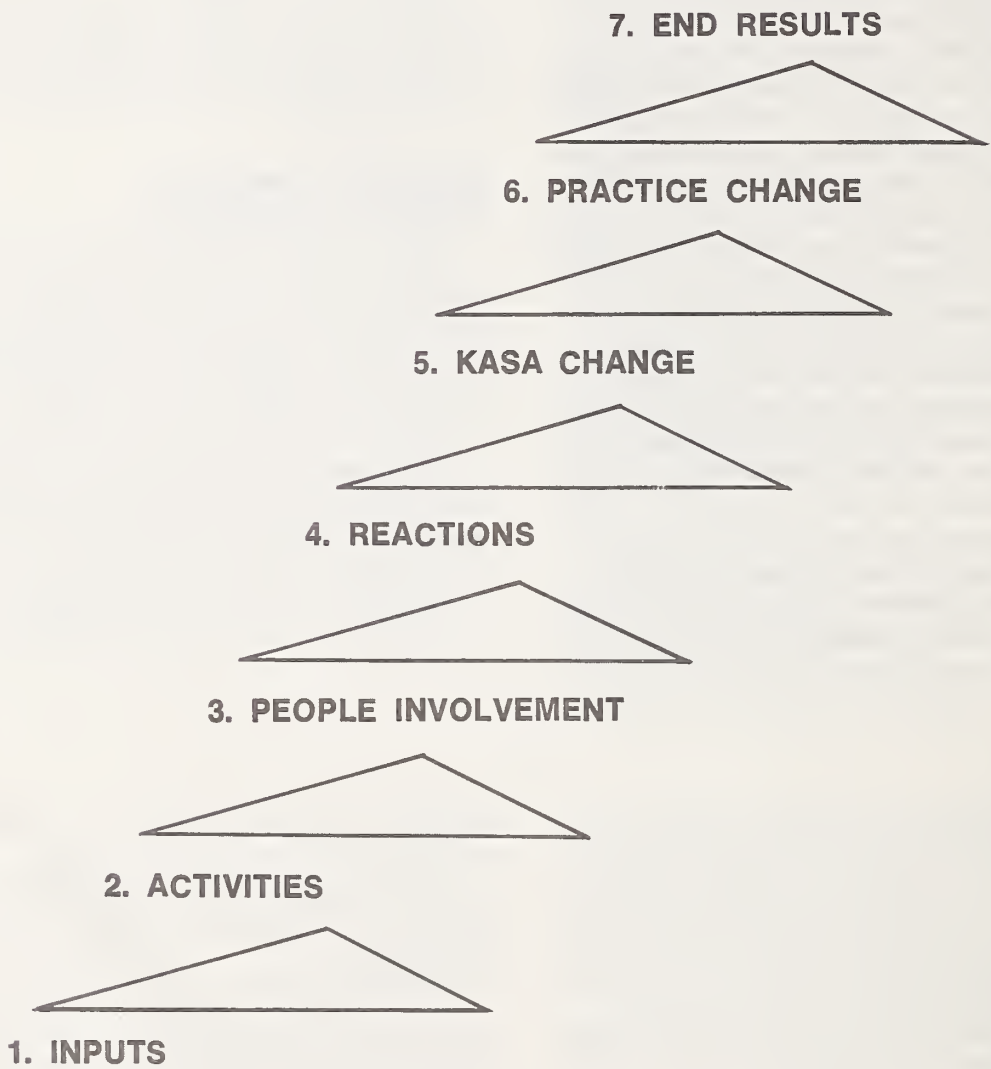
Wholey, Joseph S., John W. Scanlon, Hugh G. Duffy, James S. Fukumoto, and Leona M. Voght, *Federal Evaluation Policy: An Overview*, Washington, D.C., Urban Institute, 1970.

³ Stufflebeam, Daniel L., *op. cit.*, and Joseph S. Wholey, *op. cit.*; see also Warner, W. Keith, "Feedback in Administration," *Journal of Extension* V (Spring 1967), pp. 35-46.

⁴ Several elements of the chain have been identified by Kirkpatrick and Suchman. See Kirkpatrick, Donald L., "Evaluation of Training," *Training and Development Handbook*, Robert L. Craig and Lester R. Bittel (eds.), pp. 87-112, New York, McGraw-Hill, 1967, and Suchman, Edward A., *Evaluative Research*, New York, Russell Sage Foundation, 1967.

Figure 1.

CHAIN OF EVENTS IN EXTENSION PROGRAMS



in, and like or dislike for, the activities in which they are involved.⁵ (Reactions to activities depend on both the activities themselves and the values, learning ability, and social interrelationships of the people involved.) To the extent that participants' interest can be held, they may change their knowledge, attitudes, skills, and/or aspirations ("KASA"). Whereas attitude denotes feelings (approval or disapproval), aspiration indicates the use of feelings in goal selection or choice among alternatives.

"Practice Change" (adoption) refers to individual or collective application of acquired knowledge, attitudes, skills, and aspirations to work or life styles.⁶ But, practices are not usually adopted for their own sake; certain benefits are anticipated to accrue from individual and collective practices. Whatever benefits and consequences follow from practices may be called "End Results." These results, hopefully, include attainment of the ultimate objective(s) of Extension programs.

Before continuing, it should be acknowledged that individual or group change may not always proceed strictly in accordance with the above sequence of events. For example, reactions probably occur prior to and during participation, as well as after involvement. Also, practice change may occur before the attitude or knowledge change intended by program objectives.

A HIERARCHY FOR PROGRAM EVALUATION

In figure 2, the foregoing chain of events is converted into a hierarchy of objectives and evidence for program evaluations. Six levels of output are based upon inputs to Extension.

At each level of the hierarchy, "P" (for planned) symbolizes an objective to be reached. For example, an objective at level 3, "People Involved," could be to involve a certain number of clientele having prescribed characteristics. Placement of "P" on the sloping line at the left of each level is shown by a dot. The height of the dot opposite "P" indicates the magnitude of the objective. That is, a dot representing an objective to reach 200 clientele would be placed higher than if the objective were to reach 100 clientele. The *staircase* of objectives reaches toward solving (at the seventh level) some overall problem of clientele or the larger society. Placement of dots in figure 2 is for the sake of illustration. However, at levels 1 through 6, a basis for setting objectives is their sufficiency to move to the next higher level(s) and, finally, to the desired end results.

Figure 2 abbreviates two dimensions or broad criteria at each level; specific examples of these and other criteria are provided below:

1. At the *inputs* level, criteria are within plans (objectives) to allocate certain kinds and amounts of resources to a program, such as:

- Time of paid staff and volunteers (e.g., "five full-time equivalents per year will be expended on a consumer education program").
- Staff qualification—paid and volunteer (e.g., "all program assistants to be recruited must be 'opinion leaders'").

2. At the *activities* level, criteria are within plans to perform, through the above inputs, a certain number of specified activities in order to induce education, such as:

- Collecting and preparing educational materials (e.g., "assist volunteers in planting 20 plots to demonstrate research findings").
- Publicizing programs (e.g., "publish five newspaper notices of environmental activities").
- Transmitting subject matter content through mass media, meetings, and other events (e.g., "schedule five showings of a video tape on how to shear sheep").

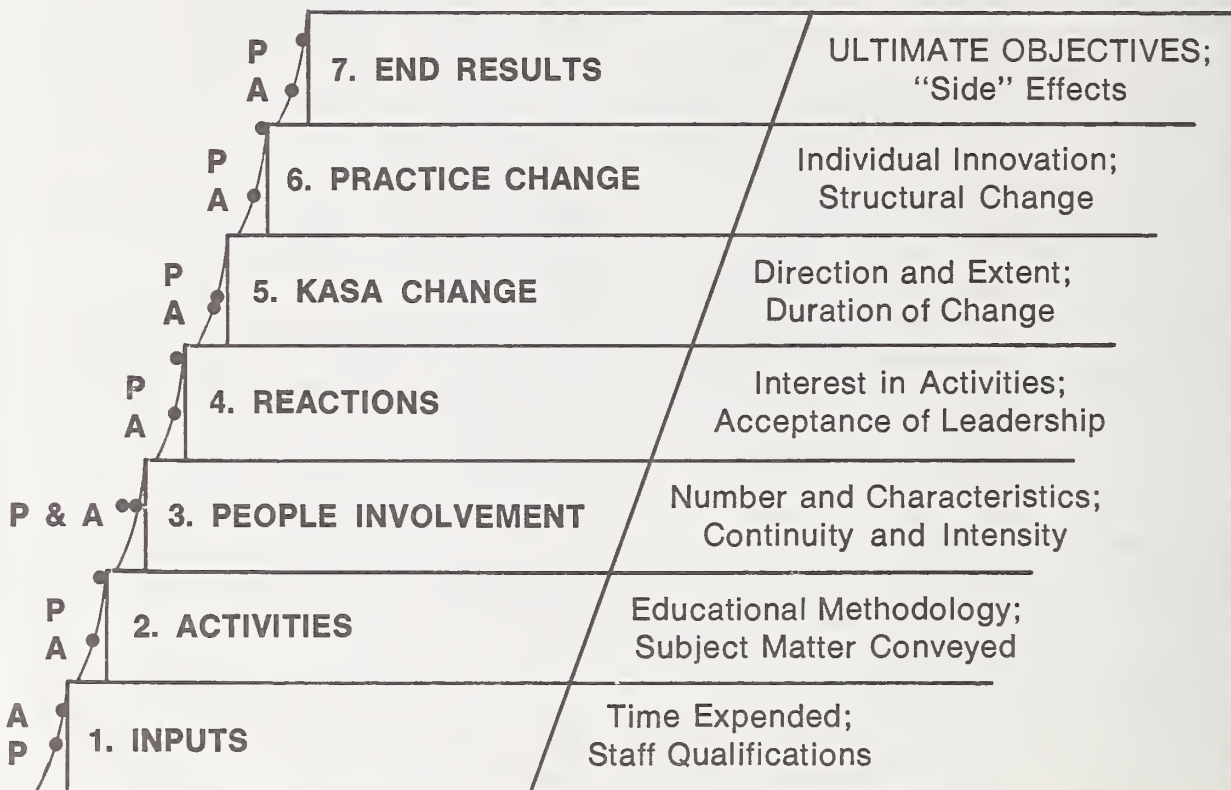
⁵ Initially, it matters little whether participants are interested in the educational content of activities; for example, some participants may at first attend and value discussions on grain production or nutrition because they enjoy the social interaction, and not because they are particularly interested in the subject under discussion.

⁶ Of course, KASA change may simply reinforce existing practices. Clientele may be expected to exhibit no practice change, as in education to prevent use of harmful drugs.

Figure 2.

HIERARCHY OF EVIDENCE FOR PROGRAM EVALUATION

Plans Compared With Achievements



KEY

P = Planned (Objective)

A = Achievement

3. At the *people involvement* level, criteria are within plans that certain types and numbers of persons, groups, or communities will be involved in the activities, such as:

- Number of participants in events, tours, meetings, or clubs (e.g., “2,000 or more 4-H club members will be enrolled in livestock projects”).
- Psychological and socioeconomic characteristics of participants (e.g., “at least 90 percent of program clientele should be low-income”).
- Continuity, frequency, and intensity of face-to-face or mediated interaction between clientele and Extension (e.g., “80 percent of new officers of community development councils should attend leadership training meetings”).

4. At the *reactions* level, criteria are within plans to obtain certain reactions to involvement in activities, in terms of:

- Interest in educational events (e.g., “there should be a minimum of 75 percent positive reactions to topics chosen for discussion at child development meetings”).
- Acceptance of persons leading activities (e.g., “leader of soybean marketing meetings should be rated as ‘highly competent’ by two-thirds of those in attendance”).

5. At the *KASA change* level, criteria are within plans that certain knowledge, attitudes, skills, and aspirations (KASA) will ensue from participants’ engagement in program activities,⁷ including:

- Direction (content) and extent of KASA change (e.g., skills—“80 percent of homemakers, rather than the present 10 percent, to be able to suitably arrange furniture in their respective homes”).

- Durability of any KASA change (e.g., knowledge—“95 percent of farmers to recall sources of safety rules for handling pesticides one year after learning about them”).

- Intensity of attitudes to be accepted (e.g., “all youth in the citizenship seminar should come to condemn very strongly the neglect to vote”).

- Height of aspiration (e.g., “each couple represented in the family resource management workshop should decide to prepare a legal will within 1 month after the close of the workshop”).

6. At the *practice change* level, criteria are within plans for certain changes in individual practices, technology, and/or social structures. These consequences of KASA change are in terms of:

- Individual innovation and adoption (e.g., “80 percent of farmers to adopt new, superior variety of wheat within 2 years of release”).

- Collective (structural) change (e.g., “25 percent of communities to establish land-use planning boards during each of 4 successive years”).

“Individual innovation” is distinguished from “structural change” in that the latter refers to change in social relationships, laws, and institutions, including associated physical facilities. For example, if a solid waste disposal system is created in a county, a *structure* within that county is changed.

7. At the *end results* level, criteria are within plans that certain effects will be achieved through practice change. These plans are called *ultimate* objectives and emphasize the prevention, checking, reduction, or solution of overall problems of:

- Individuals (e.g., “one-third of ‘isolate’ youth attending camp to gain increased self- and peer-acceptance”).

- Groups (e.g., “the community will increase to 5 percent its annual rate of real economic growth”).

REACHING PROGRAM OBJECTIVES

Figure 2 shows that actual outcomes or achievements, “A” as well as objectives, pertain to each of the seven levels discussed above. The height of the dot opposite each “A” shows the magnitude of actual outcome.

⁷ At all levels, but especially at levels 5, 6, and 7, the question of whose objectives—Extension’s or clientele’s—are involved may become an issue. The degree of consensus on objectives at these levels will depend on the adequacy of Extension program planning. See Stake, Robert E., “The Countenance of Educational Evaluation,” *Teachers College Record* 68 (April 1967), pp. 523-540.

Figure 2 also shows a variety of possible relationships between “P” and “A”. If the “A” is above the “P,” more has been attained than planned, as shown at level 1 (e.g., seven, rather than the intended five, full-time equivalents are expended on the program). If “A” is below “P,” less has been accomplished than planned, as at level 6 (e.g., over a 4-year period only 40 percent, rather than 100 percent, of communities established land-use planning boards). Of course, if plans have been exactly attained, “P” and “A” are the same, as depicted at level three (new officers of county development councils reach the objective of 80 percent average attendance at leadership training meetings).

There are many factors which enter into value judgments of programs. *However, in general, the more nearly the objectives of a program are reached, the more positive the judgment of the program, i.e., the higher the value assigned to the program.* In turn, the more a program is valued, the more likely it will be continued, intensified, or broadened (unless need for the program has been lessened due to the program’s success or to other factors).

Before continuing, it should be acknowledged that comparing objectives and achievements is by no means the only approach to evaluating Extension programs. Evaluations of program impact may be based on the entire array of program effects, whether or not related to program objectives.

“Side” effects may occur at any output level of the hierarchy, but apply especially to level 7. “Side” effects are unintentional and usually unexpected and may be beneficial or harmful. For example, new industry obtained by a community through Extension’s assistance may alter established social relationships in unexpected ways. Other approaches to program evaluation include comparing program objectives and accomplishments with the mission of Extension as an agency.⁸

SELECTION OF LEVEL OF EVIDENCE

As previously outlined, Extension programs usually have—explicitly or implicitly—objectives at several or all levels of the hierarchy depicted in figure 2. At which of the seven levels should evidence of program accomplishments be obtained in evaluating Extension’s effectiveness? Guidelines A, B, C, and D are offered to help answer this question. These guidelines, as well as others in the paper, are offered on the basis of experience and logical plausibility. Although the guidelines have not been tested systematically, they are provided in order to organize thinking about formal evaluation and to lead toward cumulation of tested principles about evaluation itself.

Guide A: Evidence of program impact becomes stronger as the hierarchy is ascended. (Of course, such evidence may indicate attainment, or lack of attainment, of objectives.) Guide A states, in effect, that evidence at the two lowest levels provides little or no measure of the extent to which clientele benefit from the program.

Level 3 merely provides one way of measuring possible opportunity for education to occur. Evidence at the “people involved” level may suggest the extent to which some kinds of benefits are being received by participants. However, evidence at this level (e.g., participation rate) does not necessarily indicate progress toward ultimate program objectives: high participation may occur for some reason unrelated to the benefits intended to accrue from the program.

Ascending to the fourth level, “reactions,” can provide somewhat better confirmation of whether given activities are helpful as intended. But such evidence indicates less satisfactorily than evidence of KASA changes the extent of progress toward ultimate program objectives. Knowledge, skills, etc., to be acquired are frequently considered as merely “stepping stones” to adoption of more desirable patterns of behavior, although there are differing philosophies on whether practice change is always necessary to successful Extension education. Practice change assessment is desirable when program objectives include patterns of: (a) utilization or application of new knowledge and skills; (b) expression of changed attitudes; and (c) follow-

⁸ Steele, Sara M., *Six Dimensions of Program Effectiveness*, Madison, Wis., Program and Staff Development, University of Wisconsin-Extension, 1972. Also see Stake, Robert E., *op. cit.*

through on new aspirations, decisions, or commitments.

Finally, assessing practice change is usually quite apart from assessing accomplishment of ultimate program objectives. Extension is often held accountable for the extent to which it is contributing to solution or checking of overall problems of clientele or the society. Therefore, *ideal* evaluation of impact of most Extension programs would probably be in terms of whether desired end results are achieved, plus assessment of any significant side effects.

However, a reason for infrequent assessment of impact at the top levels of the hierarchy is set forth in **Guide B: The difficulty and cost of obtaining evidence on program accomplishments generally increases as the hierarchy is ascended.** Evidence within lower levels of the hierarchy provides little indication of impact but is comparatively inexpensive and easily gathered. As the hierarchy is climbed, difficulty and resources required to measure actual program outcomes generally increase, due to: (a) increasingly greater difficulty in setting precise objectives as guides in obtaining accomplishment data—exclusions of alternate objectives within a level are more difficult to justify as the hierarchy is ascended; (b) increasingly scattered sources of evidence—Extension clientele often apply separately what they learn through participation in group Extension activities; (c) increasingly greater time-lag following program activities—practice changes and end results may occur months to years after Extension activities; and (d) increasing difficulty of separating Extension accomplishments from accomplishments by other sources of change—i.e., the higher in the hierarchy, the more chance that some agency, or a communication source other than Extension, had a role in bringing about any observed change.

Guides A and B both assume evidence of comparable quality from one level to another. These Guides—(A) *evidence of impact becomes stronger in ascending the hierarchy*, and (B) *more resources are required to collect evidence of accomplishment within higher levels*—are advanced only so long as the quality of evidence remains constant from level to level. The quality of the evidence is discussed later in this publication.

A PYRAMID OF EVIDENCE FOR PROGRAM EVALUATION

Figure 3 depicts a pyramid which guides toward the advantages of assessing a program at several levels of the hierarchy, including the inputs level.

Figure 3 cumulates previously discussed levels of evidence in proceeding from Evidence Clusters I to VII. Cluster I, simply the “inputs” level, constitutes an underlying component of all the other clusters of evidence. Cluster II adds a second level, “activities.” These two levels themselves constitute “building blocks” for Cluster III, and so on.

Guide C: Evaluations are strengthened by assessing Extension programs at several levels of the hierarchy including the inputs level. This guide is advanced for three reasons.

First, along with other agencies, Extension is being asked increasingly to report degree of output (levels 2 through 7) in relation to inputs or costs (level 1 of the hierarchy). This entails analysis of program delivery efficiency and of cost effectiveness or cost benefits.⁹ Clusters with higher numbers (“high” clusters) provide for analysis of program cost in relation to effectiveness criteria closer to problem solution (level 7).

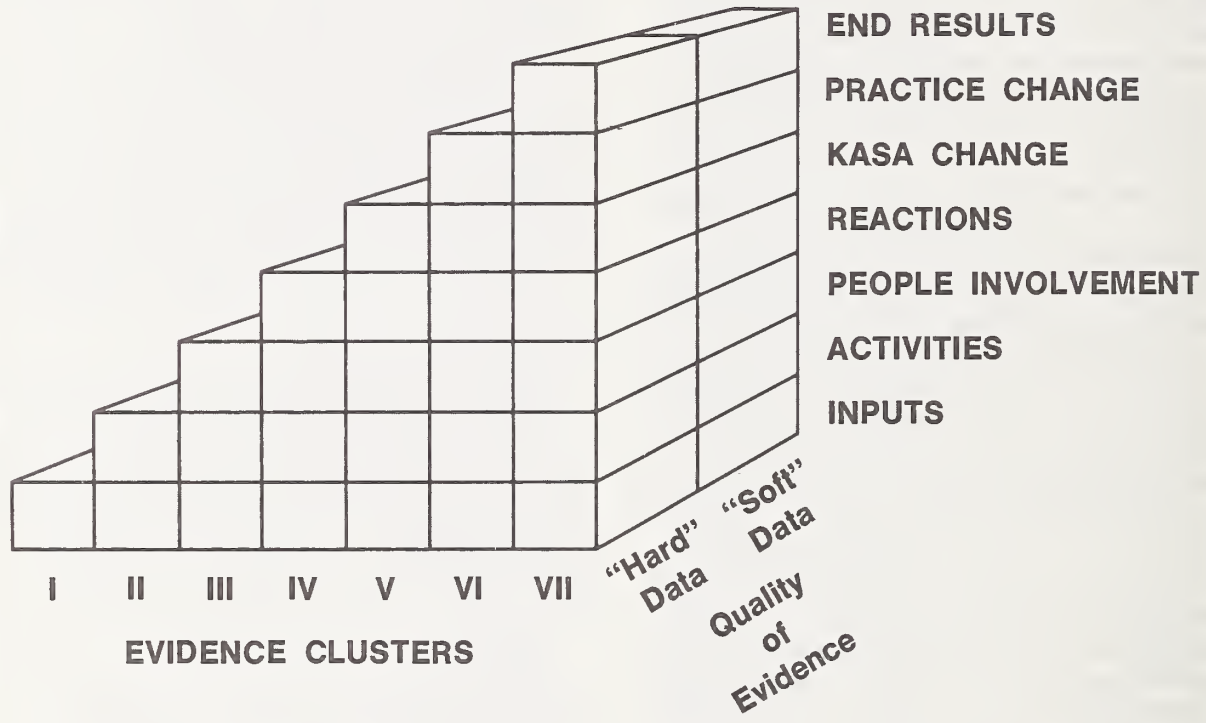
The second reason for Guide C is: the greater the number of program objectives shown to be met, including those at varying levels, the better the evidence of effectiveness. For example, evidence of intended knowledge change improves certainty that clientele adoptions of recommended practices were made for the correct reasons or because of what clientele learned through participation in Extension.

A third reason for obtaining evidence at two or more levels of the hierarchy is to check on how far the program has proceeded toward reaching its ultimate objectives. A program may fall short of inducing practice changes, but effectively induce intended KASA changes: external constraints may prevent Extension clientele from putting into practice knowledge, attitudes, skills, and aspirations acquired through participation in Extension pro-

⁹ Tripodi, Tony, Phillip Fellin, and Irwin Epstein, *Social Program Evaluation*, Itasca, Ill., F. E. Peacock, 1971.

Figure 3.

A PYRAMID OF EVIDENCE FOR PROGRAM EVALUATION



Objectives and achievements are assumed, as represented in Figure 2.

grams. Similarly, objectives for practice change may have been achieved, without sufficient time having elapsed for clientele to realize the envisioned benefits from the practice.

High clusters of evaluative evidence should be selected to the degree that resources for formal evaluation are available, as higher numbered clusters provide stronger evidence for program evaluation. In the higher clusters, one or more of the levels may be omitted, in line with the purposes or constraints of compiling evidence for evaluation.

The paper to this point may be partly summarized and also related explicitly to the chief purpose of program evaluation, by stating **Guide D: The higher the cluster of evidence for program evaluation, the more useful the evidence for making decisions on present and future programming.**

EVALUATION CRITERIA AND QUALITY OF EVIDENCE

Assessing program effectiveness generally requires specific criteria which can provide a basis for *measuring* the extent to which program objectives have been attained.

Criteria within program objectives are generally definitions or subdivisions of objectives at each level of the hierarchy.¹⁰ Criteria are a primary basis for selection of evidence as to the extent of accomplishment of objectives. For example, if the *ultimate aim* of a program (level 7 of the hierarchy) is to achieve "desirable land-use," how would "desirable land-use" be defined? Would it be defined in terms of trade-offs among preferred (a) "living space," (b) "population growth," (c) "economic growth," and (d) "environmental status"? If so, how would (a), (b), (c), and (d) be defined? Repeated subdivision of the components of (a), (b), (c), and (d) would continue until criteria are sufficiently specific and clear to guide the selection of adequate evidence on the extent to which ultimate aims of the program have been achieved.

The process of defining specific criteria for

evaluation is essentially one of moving from broad to specific objectives at each level of the hierarchy. Therefore, planning for obtaining evaluative evidence can and should occur simultaneously with the process of preparing multiyear programs, annual plans of work, and learning activities.

Guide E: Evaluation is strengthened to the extent the specific criteria for evaluation are defined prior to conduct of the Extension program. Specific criteria are needed in order to obtain quality evidence on degree of attainment of program objectives: (a) prior to program activities ("benchmark" evidence), and, (b) following such activities. Early timing in planning for evaluation can clarify program objectives and, thus, also strengthen the planning and conduct of Extension programs. Timing in obtaining evidence will be discussed in some detail later.

Evidence on the extent of accomplishment of objectives may vary in quality. Variation in quality of evidence is often referred to as "hard" versus "soft" data. Data (i.e., observations) are "hard" to the extent that they are valid, representative, and quantified.¹¹ Figure 3 indicates that "soft" or "hard" data (or both) may be collected at each level of the hierarchy. It should be emphasized that "hard" and "soft" data constitute a continuum; a dichotomy is depicted simply for the sake of convenience.

The degree of "hardness" of data actually selected depends upon trade-offs between ideal data for the evaluative purpose at hand and the resources available. Hard data are usually ideal; however, "hard" data are also more expensive and difficult to obtain and should be collected only when the benefits to decisionmaking anticipated from superior evidence clearly outweigh the costs of obtaining such evidence.

There are many situations where "soft" data on degree of accomplishment of objectives are all that can be obtained; for example, program participants, and especially nonparticipants, are often unwilling

¹⁰ Criteria for evaluating program impact may be unrelated to program objectives. For example, criteria may be based on philosophical, ethical, or personal considerations.

¹¹ For an introduction to validity, quantification, and representativeness, see Selltitz, Claire, Marie Jahoda, Morton Deutsch, and Stuart W. Cook, *Research Methods in Social Relations*, New York, Holt, Rinehart and Winston, 1961.

or unable to be observed or to respond to instruments which require detailed answers and extensive time for completion.

Figure 4 illustrates the three principal dimensions of “hard” versus “soft” data. Observations are *valid* to the extent that they truly reflect the characteristics of individuals, groups, or situations under study. For example, regarding the measurement of knowledge (level 5), validity of responses by workshop participants to the following questionnaire item would be rather uncertain: “Please indicate whether you can recognize potassium deficiency in wheat plants: (1) _____ ‘very confident I can,’ (2) _____ ‘fairly confident I can,’ (3) _____ ‘not sure I can.’” (A participant’s selection of one of the three responses could be observed and, therefore, be considered as data). A wholly valid measure of participants’ actual knowledge would entail direct observation of the degree to which they can, in fact, accurately identify potassium deficiency under given conditions, such as developmental stage of plants presented, etc. Precise definitions would be needed to specify observable actions indicating correct recognition of potassium deficiency.¹²

Nonvalidity of data may arise from several sources. For example, awareness of program evaluation by participants may cause them to speak or act as they think they are expected to for the sake of the evaluation, rather than in accordance with their own inclinations. Lack of validity may also arise from faulty instruments of observation, from observing too small a range of actions by Extension participants, and from perceiving participants’ actions inaccurately due to personal bias.

Validity of observations is demonstrated by the extent to which they are consistent with other relevant evidence concerning characteristics of individuals, groups, or situations under study. **Guide F: Evaluations are strengthened to the extent that validity of observations has been demonstrated.**

Although true differences in characteristics of

units may be *observed* as differences (validity), the question as to amount or degree of difference remains. This poses the dimension of *quantification*. The degree of difference may be shown by the assignment of numerals to represent quantities. Thus, quantitative data indicate *how much* difference there is in individuals or structures which are observed.

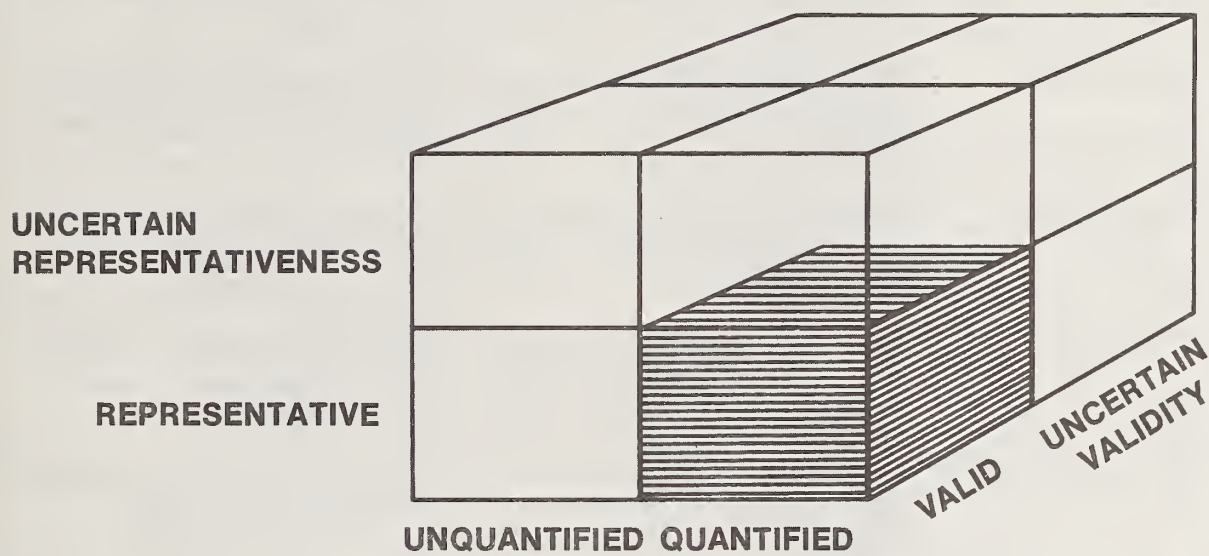
Suppose that participants of a tour of feed lots are asked to rate various displayed practices to minimize environmental pollution as “acceptable” or “unacceptable” for their own use (a measure of their attitude toward each practice). Such a rating does not permit measurement of whether one practice acceptable to the individual exceeds his acceptance of another. The participant’s responses could be quantified by asking him to rate each pollution control practice on a scale of “zero through 10”. “Zero” could represent “totally unacceptable,” and “10,” “totally acceptable,” with varying degrees of acceptability represented by numbers 1 through 9.

The third aspect of hard data is that of *representativeness*. Representativeness is the extent to which observations concerning individuals, groups, or situations under study apply to some total population of individuals, groups, or situations. Representativeness may be obtained by conducting a census or selecting a representative sample. In program evaluation, a census obtains information from (or on) all the actual or potential program participants. A representative sample may be chosen so that the information obtained corresponds closely enough, for the purposes at hand, to comparable census findings. Every tenth recipient of a consumer economics newsletter might be a sufficiently representative sample for the purpose of evaluating the newsletter.

The weight given to an evaluation in making a program decision should depend upon hardness of the evidence. **Guide G: The harder the evidence for evaluation, the more an evaluation may be relied upon in program decisionmaking.** Table 1 shows examples of “hard” and “soft” data at each level of the hierarchy.

¹² Mager, Robert A., *Preparing Instructional Objectives*, San Francisco, Fearon Publishers, 1962.

Figure 4.
QUALITY OF EVIDENCE
"Hard" Versus "Soft" Data



KEY



-  "Hard" Data
-  "Soft" Data

Table 1
Examples of “Hard” and “Soft” Data in a Hierarchy
of Evidence for Program Evaluation

	Examples	
	“Hard” data	“Soft” data
7. End results	Trends in profit-loss statements, life expectancies, and pollution indexes	Casual perceptions of changes in quality of health, economy, and environment
6. Practice change	Direct observation of use of recommended farm practices over a series of years	Retrospective reports by farmers of their use of recommended farm practices
5. KASA change	Changes in scores on validated measures of knowledge, attitudes, skills, and aspirations	Opinions on extent of change in participants’ knowledge, attitudes, skills, and aspirations
4. Reactions	Extent to which random sample of viewers can be distracted from watching a demonstration	Recording the views of only those who volunteer to express feelings about demonstration
3. People involvement	Use of social participation scales based on recorded observations of attendance, holding of leadership positions, etc.	Casual observation of attendance and leadership by participants
2. Activities	Pre-structured observation of activities and social processes through participant observation, use of video and audio tapes, etc.	Staff recall of how activities were conducted and the extent to which they were completed
1. Inputs	Special observation of staff time expenditures, as in “time and motion” study	Staff’s subjective reports regarding time allocation

PROXY MEASURES

Extension frequently lacks sufficient resources to obtain quality evidence of its extent of effectiveness, especially at higher levels of the hierarchy. In such cases, *inferences* of the degree to which objectives are attained can be made if *proxy* or substitute measures have been established.¹³ Proxy measures are based on research-tested relationships between the achievement of objectives at lower and higher levels of the hierarchy, e.g., between KASA change and desired practice change in a youth community development program. On the basis of such previous research, reaching a lower level objective in a program permits inferring or predicting attainment of a higher level objective. Of course, caution *must* be exercised as to how far previous research can be generalized as a basis for assessing program effectiveness.

With their more confined scope and variation, demonstration and pilot projects permit collection of "high" evidence clusters with resources comparable to those necessary for collection of "low" evidence clusters on full-scale programs. An efficient strategy for an agenda of formal evaluation is this: collect high clusters on pilot projects and, in so doing, identify within lower levels of the hierarchy proxy measures of impact. These proxy measures can provide a basis for interpretation of subsequent low evidence clusters collected on any ensuing full-scale program. Similarly, if Extension can evaluate full-scale programs periodically through high cluster evaluations, then, between such evaluations, low clusters can be used to make inferences about achievement of objectives at higher levels of the hierarchy.

Through application of the above strategies, a schedule of evaluations can be designed to provide systematically over a cycle of years for efficient formal evaluation of Extension's programs or program components. ***Guide H: The efficiency of program evaluation can be increased through studies which identify proxy measures.***

DESIGNS FOR IDENTIFYING SOURCE OF IMPACT

Study designs suggest schemes for collecting evidence of Extension's impact. Designs vary in strength of scientific evidence regarding the extent to which KASA change, practice change, or end results were brought about through Extension rather than through other sources of change. Of course, Extension often works along with other agencies and institutions in addressing problems.

There is scientific evidence of Extension's impact, to the degree that evidence can exclude or take into account other possible causes of achievement of program objectives (e.g., other programs, chance events, maturation of participants, effects of being observed or tested before the program, special motivation of clientele involved in Extension, etc.).¹⁴

Guide I: A study's usefulness for program decisionmaking is enhanced to the extent that it can identify Extension's degree of contribution to achievement of program objectives.

The following are only a few of the possible study designs. First presented is the field experiment, which provides strongest scientific evidence of the degree to which observed change is produced through Extension. Other designs are presented in order of their capability of identifying the degree to which Extension contributes to observed attainment of program objectives. The designs are not necessarily limited to the way in which they are described below: each may be more or less complex in being adapted to varying conditions. The designs are defined and illustrated below to show a range of possibilities for identifying Extension's contribution to change. Finally, the designs are described in relation to program objectives, in order to show their relevance to program evaluation as defined in this publication.

¹³ Wholey, Joseph S., John W. Scanlon, Hugh G. Duffy, James S. Fukumoto, and Leona M. Voght, *Federal Evaluation Policy: An Overview*, Washington, D.C., Urban Institute, 1970.

¹⁴ Campbell, Donald T., and Julian C. Stanley, *Experimental and Quasi-Experimental Designs for Research*, Chicago, Rand McNally, 1963. Also Stouffer, Samuel A., "Some Observations on Study Design," *The American Journal of Sociology* 55 (January 1950), pp. 356-359.

16 THE FIELD EXPERIMENT

Two instances of field experimental evaluation studies in Extension are: (a) the national impact evaluation of *Mulligan Stew*, a televised nutrition program for youth,¹⁵ and (b) a study of Extension training impact on managers of Iowa retail farm supply firms.¹⁶

The field experiment requires making the program available to clientele selected randomly (through chance alone) from some audience. The part of the audience selected for no exposure to the program is the "control group." For example, farm and family Extension aides could be assigned to disadvantaged rural residents in half of the counties of a State. These counties, selected at random, would contain the program group of disadvantaged rural residents. The other counties would contain the control group. Observations before and after the program activities within both the program and control groups are usually desirable in field experiments. However, observations only *after* the activities are permissible in the conduct of field experiments and may be preferred under some circumstances.

Figure 5 depicts possible observations in a field experiment.¹⁷ In figure 5, levels of the hierarchy in which no observations are made are represented by broken lines. Observations prior to program activities ("before" observations) are made simultaneously at levels 5, 6, and 7 in both the program and control groups. The "situation or benchmark" in each group is the same, as shown by the identical location of "A_b" relative to the sloping lines of levels 5, 6, and 7. Turning now to "during observations," the action strategy and reception actually occurred as planned, as shown by the coincidence of "P" and "A" at levels 1, 2, 3, and 4 in the program group. "After" observations of KASA, practices, and end

results are made as soon as it is reasonable to expect that the intended changes at these three levels have occurred.

The interrelationships among objectives and observed achievements shown in figure 5 suggest an effective Extension program. First, "A" reaches "P" in the four top levels of the *program* group. Secondly, although each "A" in the control group is higher than in the "before" situation, the rise is less than the rise of the corresponding *program* group "A." The contribution of sources of change other than Extension is shown by comparing the "before-after" observation within the control group (A_b compared with A). A "significance" test can gauge the odds that any greater increase in program group achievement over that of the control group was brought about by the presence of the program rather than by uncontrolled factors or chance.

The field experiment should be used when it is essential to have maximum certainty about the extent of Extension program impact. In many situations, the field experiment is unattainable because of complexity or cost, or undesirable because of ethical or political considerations. Under such conditions, it is necessary to settle for designs which provide evidence less conclusive of Extension's impact.

MATCHED SET DESIGN

The comparison set design is similar to the field experiment except that program availability to a portion of the potential audience is on other than a *random* basis. Rather, a program group (set) and a comparison set are usually selected on the basis of their similarity. For example, (a) a study in New York State¹⁸ compared progress of farmers in an Extension farm management program and progress of similar

¹⁵ Shapiro, Sydelle S., Richard L. Bale, Vince Scardino, and Tom Cerva, *An Evaluation of the Mulligan Stew 4-H Television Series for Extension Service, USDA*, Vols. I, II, III, and IV, Cambridge, Mass., Abt Associates, 1974.

¹⁶ Warren, Richard, George M. Beal, and Joe M. Bohlen, *The Experimental Dealer Training Program*, Ames, Iowa State University, Rural Sociology Report 56, 1966.

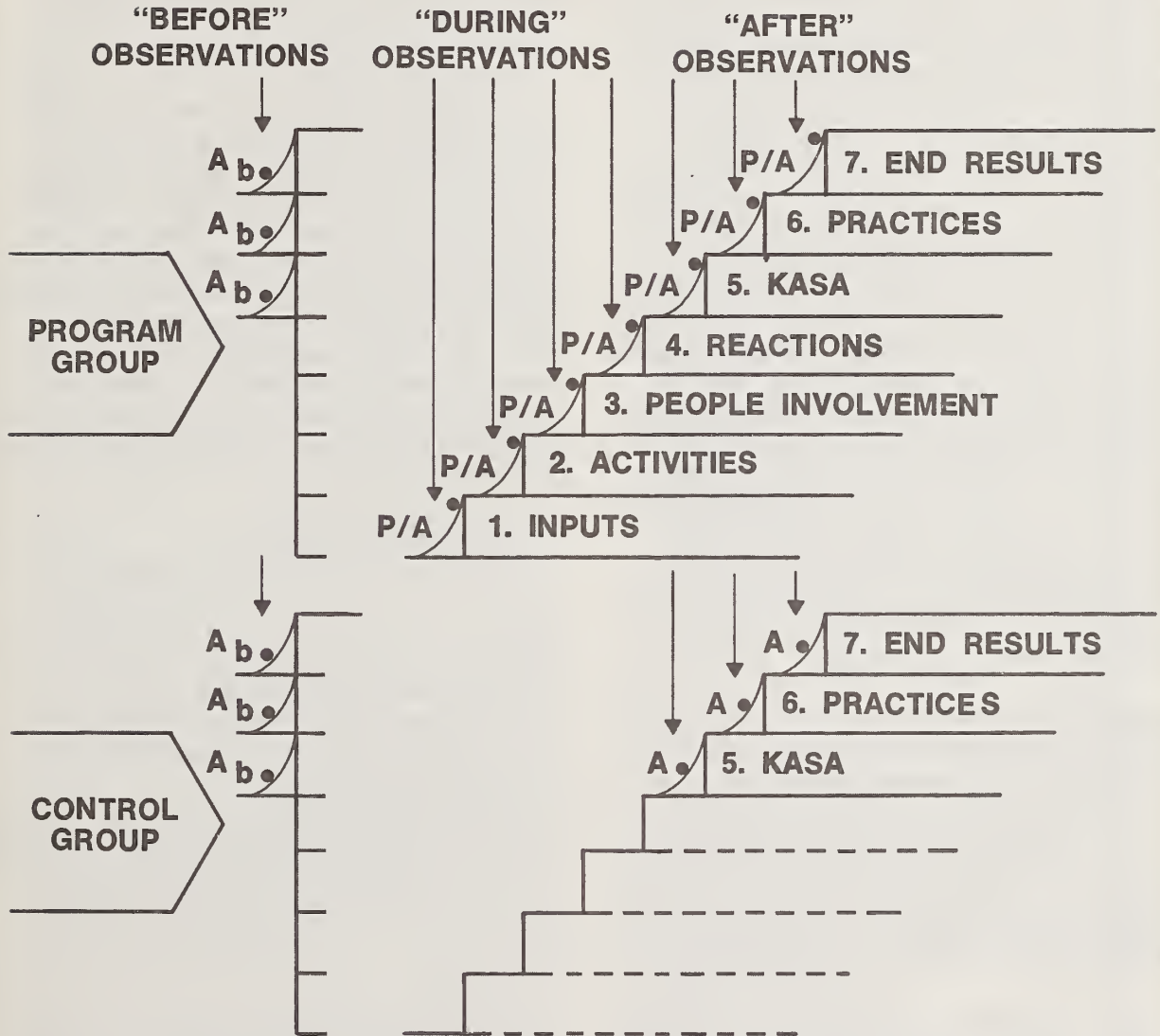
¹⁷ Bennett, Claude F., and Robert C. Leonard, "Field

Experimentation in Rural Sociology," *Rural Sociology* 35 (March 1970), pp. 69-76.

¹⁸ Alexander, Frank D., and James W. Longest, *Evaluation of the Farm Management Phase of the Farm and Home Management Program in New York State*, Ithaca, Office of Extension Studies, New York State Extension Service, 1962.

Figure 5.

POSSIBLE OBSERVATIONS IN A FIELD EXPERIMENT



KEY

P = Planned

A = Actual

A_b = Actual before program

NOTE: "During" includes steps 1-4;
 "After" includes steps 5-7.

18 nonprogram farmers; (b) a Maryland study¹⁹ compared progress of families in an intensive nutrition education program with the progress of families who were friends of program families, but who were not themselves program participants; and (c) the impact of a community improvement program in Kentucky was studied through matching program and nonprogram communities.²⁰

The basic limitation of the matched set design is that matching can be only partial, and not complete.²¹ To the extent that matching is incomplete, the matched set design fails to identify accurately Extension's contribution to change (as compared to other sources of change). The matched set design does not provide for statistical tests to determine the odds that extraneous factors are responsible for any greater change in program set "A" than in comparison set "A." Statistical techniques such as co-variance analysis or multiple regression can correct partially for such extraneous factors, but cannot substitute fully for the random assignment in the field experiment. The matched set design should seek to identify factors in addition to Extension which may effect change, so that at least these factors may be accounted for statistically in assessing Extension's degree of contribution to accomplishment of program objectives.

TIME-TREND STUDIES

These studies follow clientele's KASA change, practice change, or problem solution over an extended (e.g., multi-year) period. There are two major variations of this method. The first is time-trend projection of preprogram data vs. actual observations after program implementation. Program impact is identified as the difference between observed "after" program conditions

and projected conditions based on rates of change from time periods prior to the program.²² Of course, to account for the amount of change which has occurred, it is necessary to look for plausible explanations other than the Extension program. This design is appropriate when there is a trend that seems likely to have continued if the program had not been introduced (e.g., rate of increase in average number of pounds of milk per cow per year).

A second type of time-trend study is one which obtains repeated measurement of clientele progress relative to program objectives. A prime example of this design is the national evaluation of the Expanded Food and Nutrition Education Program. In this program, clients are enrolled as program participants, thus facilitating observation of their KASA change, practice change, and degree of problem solution over the length of their participation. Observation of self-reported food consumption has been made beginning with entry of the client into the program and every 6 months thereafter.²³

"BEFORE-AFTER" STUDY

This design requires observations both before and after an Extension program, as could be shown by the program group portion of figure 5; no comparison set or "control group" is used. The before-after design has been used in many Extension studies and is well exemplified by the evaluation of a Texas Extension program for low-income farmers.²⁴

The "before-after" design tests only partially the extent to which any changes at higher levels in the hierarchy are produced by Extension inputs, activities, etc. But, it is plausible that Extension produced part of any observed impact, to the

¹⁹ Green, Lawrence W., Virginia Li Wang, and Paul H. Ephross, "A Three-Year Longitudinal Study of the Impact of Nutrition Aides on the Knowledge, Attitudes and Practices of Rural Poor Homemakers." Paper presented to American Public Health Association, Atlantic City, November 1972.

²⁰ Street, Paul, *The Appalachian Community Impact Project: Comparison of Change Among the Adults and Youth*, Lexington, Ky., Cooperative Extension Service, 1972.

²¹ Alexander, Frank D., "A Critique of Evaluation," *Journal of Extension* 3 (Winter 1965), pp. 205-212.

²² Hatry, Harry P., Richard E. Winnie, and Donald M. Fisk, *Practical Program Evaluation for State and Local Government Officials*, Washington, D.C., The Urban Institute, 1973.

²³ Economic Research Service, *The Expanded Food and Nutrition Education Program 1969-1973*, Washington, D.C., U.S. Department of Agriculture, 1975.

²⁴ Ladewig, Howard, and Vance W. Edmonson, *The Effectiveness of Nonprofessionals in Cooperative Extension Education for Low-Income Farmers*, College Station, Tex., Texas A&M University, 1972.

degree that other possible source of KASA change, practice change, etc., may be ruled out logically. However, simple comparison of “before” and “after” program data may be misleading due to unusual or normal fluctuations such as seasonal variations.

Although the designs described above—(a) field experiment, (b) matched set design, (c) time-trend studies, and (d) before-after study—are desirable, implementation of these designs can be cumbersome, expensive, and difficult to complete soon enough to assist decisionmaking on future programming.

The general use of data “on the hard side” in the four designs above accounts for much of their expense and time consumption. Moreover, as the designs above select evidence increasingly high in the hierarchy: (a) the longer it is usually necessary to wait “till the data are in;” (b) the more things can complicate the study, such as attrition of program participants; and (c) the more expensive the study is to complete. A frequent upshot is use of the designs below which are less capable of controlling for “rival explanations” (i.e., attributing observed changes to sources other than Extension).

THE SURVEY

In comparison with experimental, matched-set, time-trend, and “before-after” designs, the survey design requires fewer resources per program participant observed. No “before” observations are made in the survey, which may be depicted by the “after” (and also perhaps by the “during”) observations shown in figure 5. Surveys in program evaluation may compare Extension clientele and nonclientele within higher levels of the hierarchy.

Or, the survey may compare at one point in time achievement of program objectives by Extension clientele with different characteristics, including

high versus low degree of program involvement. In such a survey, participants with a low degree of involvement constitute a partial substitute for a “comparison set.”

Primarily because of lack of situational data *prior* to an Extension program, the survey generally provides rather weak conclusions about the extent to which Extension, rather than other forces, produces any observed differences between Extension clientele and nonclientele. Limitations to such inferences from surveys include self-selection as a participant in Extension and the effect on survey observations of any “drop-outs” from the Extension program. Even with complex statistical analysis, the survey usually provides limited capacity to account for the degree to which Extension produces achievement of higher level objectives.

An important use of the survey is to collect data on *perceptions* or *opinions* about the activities and outcomes of Extension programs.²⁵ A random sample of opinions as to effectiveness of Extension programs may be evidence sufficient to meet evaluative needs.

Opinions may be obtained regarding a wide variety of areas, such as: (a) the extent to which Extension program objectives have been achieved; (b) the extent to which Extension and other actors, agencies, etc., have produced given outcomes; and (c) the degree to which Extension clientele are satisfied with Extension’s programs. Numerous studies using this methodology have been conducted.²⁶ A modified form of survey is elicitation of retrospective reports on participants’ status prior to their program participation. These reports provide a partial substitute for “before” measurements.²⁷ Retrospective reports are generally less reliable than responses reflecting the present, except where substantiating records are available.

Despite its many limitations, the survey design

²⁵ Hays, Samuels P., Jr., *Evaluating Development Projects*, Paris, Imprimerie Boudin, United Nations Educational, Scientific and Cultural Organization, 1965.

²⁶ Rose, Donald W., *A Comparative Study of Two Patterns of Cooperative Extension Organization in Colorado and Their Association with Goal Achievement, Job Satisfaction and Clientele Satisfaction*, Ph. D. Dissertation, University of Utah, 1971.

Davie, Lynn, Terry Patterson, Dorothy MacKerachey, and Richard Cawley, *SHAPES: Shared Process Evaluation System*, Toronto, Can., Ontario Institute for Studies in Education, 1975.

²⁷ Oldham, Marvin D., and Claude F. Bennett, *A Concerted Effort in Rural Development: Analysis and Evaluation*, Stillwater, Cooperative Extension Service, Oklahoma State University, 1975.

lends itself to many evaluation situations. It is comparatively simple and flexible. The survey may be employed after a program is implemented, without the prior evaluative planning required by some other designs. Finally, compared in terms of the number of persons or participants included in a study, opinion surveys are usually less expensive than are the study designs discussed previously.

THE CASE STUDY

Case studies observe intensively one or only a few selected individuals, groups, or communities. Observation may involve examination of existing records, interviewing, or participant observation.

Case studies seldom carry the rigor or formality of the preceding designs. They often use soft data (especially in terms of questionable representativeness) and seldom employ statistical analysis. In contrast to the designs discussed above, few, if any, explicit comparisons are made: the case selected for study is compared only implicitly with other cases casually observed or remembered.

The weakest form of the case study, as used in Extension evaluation, is the isolated "success story," which documents the progress of only one or several clientele. Such case studies provide weak scientific evidence of Extension's impact in a community, state, or nation, because: (a) even if data on each case is valid, the cases may not be representative of Extension clientele, and (b) the question of how much progress clientele and potential clientele would probably have made without Extension's aid is usually not answered satisfactorily.

Stronger case studies are those conducted by

outside observers using their own perceptions of program process and impact and drawing on the observations of key observers.²⁸

The case study can draw together many diverse pieces of information into a unified interpretation and may provide important evaluative insights. Thus, the case study can provide leads regarding the conduct and interpretation of studies which use more definitive designs.

Table 2 summarizes major characteristics of the six designs above.

USING AND APPRAISING EVALUATION STUDIES

Evaluations of program effectiveness are utilized most fully if their implications for decisionmaking are noted explicitly. *Guide J: Usefulness of evaluation reports is maximized when they include alternatives and recommendations for future program development.* Interpretation of evaluation findings for decisionmaking should include appraisals of the quality and completeness of the evaluation study.

The collection, analysis, and use of evidence in judging degree of program effectiveness should itself be assessed for effectiveness. If acquisition and use of evidence on program impact is viewed as an "activity" through "inputs," then a number of questions follow, based on the hierarchy for evaluation presented in this paper. Examples of these questions are: "What has been learned by

²⁸ Niederfrank, E. J., Francis S. Mansue, and Chester R. Smith, *Helping New Jersey Urban Youth Help Themselves*, New Brunswick, N.J., Cooperative Extension Service, Rutgers University.

**Table 2:
Characteristics of Designs for Analyzing
Impacts of Extension Programs**

Evaluation design	Observations			Comparison set		Evidence can apply broadly
	"Before"	"During"	"After"	Used	Randomly assigned	
Field Experiment	Maybe	Yes	Yes	Yes	Yes	Yes
Matched set	Yes	Yes	Yes	Yes	No	Yes
Time trend	Yes	Yes	Yes	No	—	Yes
"Before-After"	Yes	Yes	Yes	No	—	Yes
Survey	No	Maybe	Yes	Maybe	—	Yes
Case study	Maybe	Maybe	Yes	No	—	Maybe

the collection and analysis of data, in relation to degree of expected improvement in knowledge about program effectiveness?" "Have program decisions been influenced by knowledge of program effectiveness acquired through evaluation studies?"

Appraisals of evaluation studies can suggest needs for further program evaluation, or related analyses, to assist in specific decision issues.

SUMMARY AND CONCLUSIONS

The major purpose of evaluations is to assist in program decisions. Formal evaluations are worth doing only if they have a chance of affecting such decisions.

This publication presents options and guidelines relative to: (1) selection of strength of evidence of Extension's impact and (2) resources required for obtaining evidence. Selection of strength and expense of evidence on program effectiveness vary with informational needs and resources of decision-makers.

Selection of evidence for evaluation studies should be guided by the following questions:²⁹

1. Which *levels* of evidence for program evaluation are desired for decisionmaking relative to program continuation, direction, size, methodology, audience, etc.?
2. How "*hard*" should the evidence be, and what kind of *study design* is needed to assist materially in decisionmaking?
3. Are resources available to obtain desired level(s) and hardness of data, and to implement the desired study design?
4. If the answer to question 3 is "yes," then fine!
But if the answer is "no," then:
 - a. Can additional resources be obtained to acquire the needed evidence? If the answer is again "no," then:
 - b. Can decisionmakers use evidence from a lower level, softer evidence, or evidence from a weaker study design?

Adequate judgments of program value and sound program planning decisions can be made only by comparing clear criteria and sufficient evidence regarding program accomplishments.

ACKNOWLEDGMENTS

Gratitude is expressed to the following persons for criticisms and suggestions in preparation of this publication:

Patrick Borich, State Leader, Extension Research and Education, Minnesota Agricultural Extension Service, University of Minnesota
John Fedkiw, Assistant Director, Office of Management and Finance, U.S. Department of Agriculture
Marilyn A. Jarvis-Eckert, former State Program Leader, Nutrition Education, Center for Extension and Continuing Education, West Virginia University
Opal Mann, Assistant Administrator, Home Economics, Extension Service, U.S. Department of Agriculture
Sara M. Steele, Professor of Agricultural and Extension Education, University of Wisconsin
W. Keith Warner, Professor of Sociology, Brigham Young University.

The author also appreciates assistance from the following:

Harold J. Alford, Educational Testing Service, Princeton, N.J.
Charles Beer, Extension Service, USDA
Milton Boyce, Extension Service, USDA
Patrick G. Boyle, CES*, University of Wisconsin
Damaris Bradish, CES*, University of Arizona
Mary L. Collings, Extension Service, USDA (Retired)
Robert Dotson, Agricultural Extension Service, University of Tennessee
Robert Frye, Economic Research Service, USDA
J. J. Lancaster, Extension Education, University of Georgia
Robert C. Leonard, Sociology, University of Arizona
James J. McAllister, CES*, Oregon State University
David J. Miller, CES*, North Dakota State University
Robert W. Miller, Office of Research and Development, West Virginia University
Clyde Richardson, CES*, Colorado State University
Richard M. Ryckman, Psychology, University of Maine
Douglas Sjogren, Human Factors Research Laboratory, Colorado State University
Hoyt Warren, CES*, Auburn University (Retired)
Maria I. Cabrera, Lorraine May, June Sullivan, and Almira Swygart, Extension Service, USDA.

²⁹ Bennett, Claude F., "How to Analyze Impacts of Extension Programs," Program and Staff Development, Extension Service, USDA, 1974.

* Cooperative Extension Service

