

Evaluation in Practice

A Methodological Approach

edited by

Richard D. Bingham

Cleveland State University

Claire L. Felbinger

American University

SECOND EDITION

Evaluation in Practice

A Methodological Approach

SECOND EDITION

Richard D. Bingham

Cleveland State University

Claire L. Felbinger

American University

CHATHAM HOUSE PUBLISHERS

SEVEN BRIDGES PRESS, LLC

NEW YORK · LONDON

Experimental Designs

THERE ARE EXPERIMENTAL DESIGNS and there are experimental designs. A distinction must be made between experimental designs and the quasi-experimental designs that are discussed in Part III. The concern here is with the most powerful and “truly scientific” evaluation design—the controlled, randomized experiment. Essentially three “true” experimental designs can be found in the literature:

- (1) the pretest-posttest control group design,
- (2) the Solomon Four-Group Design, and
- (3) the posttest-only control group design.

Actually, a fourth variation—the factorial design—was discussed in chapter 2.

As Peter Rossi and Howard Freeman (1985, 263) note, randomized experiments (those in which participants in the experiment are selected for participation strictly by chance) are the “flagships” in the field of program evaluation because they allow program personnel to reach conclusions about program impact (or lack of impact) with a high degree of certainty. These evaluations have much in common with experiments in the physical and biological sciences, particularly as they enable the research results to establish

causal effects. The findings of randomized experiments are treated with considerable respect by policymakers, program staff, and knowledgeable publics.

The key is randomization, that is, random assignment. True experimental designs always assign subjects to treatment randomly. As long as the number of subjects is sufficiently large, random assignment more or less guarantees that the characteristics of the subjects in the experimental and control groups are statistically equivalent.

As David Nachmias points out, the classical evaluation design consists of four essential features: comparison, manipulation, control, and generalizability (1979, 23–29). To assess the impact of a policy, some form of comparison must be made. Either a comparison is made of an experimental group with a control group or the experimental group is compared with itself or with some selected group before and after treatment. In a true experimental design, the experimental group is compared to a control group.

The second feature of an evaluation design is manipulation. The idea is that if a program or policy is actually effective, the individuals (or cities, or organizations) should change over the time of participation. If we

are able to hold all other factors in the world constant during the evaluation, then the change in policy (manipulation) should cause a change in the target exposed to it (individuals, cities, or organizations).

The third feature of the experimental design—control—requires that other factors be ruled out as explanations of the observed relationship between a policy and its target. These other factors are the well-known sources of internal invalidity discussed in chapter 2: history, maturation, testing, instrumentation, statistical regression, selection, and experimental mortality. As Nachmias points out, these sources of internal invalidity are controlled through randomization (1979, 27–28).

The final essential feature of the classical design is generalizability, or the extent to

which research findings can be generalized to larger populations and in different settings. Unfortunately, the mere use of the controlled, randomized experimental design will not in itself control for sources of external invalidity or the lack of generalizability.

The three chapters in this section examine and critique three studies that illustrate the use of the pretest-posttest control group design, the Solomon Four-Group Design, and the posttest-only control group design.

References

- Nachmias, David. 1979. *Public Policy Evaluation: Approaches and Methods*. New York: St. Martin's.
- Rossi, Peter H., and Howard E. Freeman. 1985. *Evaluation: A Systematic Approach*. 3d ed. Beverly Hills, Calif.: Sage.

Pretest-Posttest Control Group Design

THE PRETEST-POSTTEST control group experiment, the classical experimental design, consists of two comparable groups: an experimental group and a control group. Although a number of authors use the terms “control group” and “comparison group” interchangeably, in a strict sense they are not. True control groups are formed by the process of random assignment. Comparison groups are matched to be comparable in important respects to the experimental group. In this book, the distinction between control groups and comparison groups is strictly maintained.

When the pretest-posttest control group design is used, individuals are randomly assigned to one of the two groups. Random assignment of members of a target population to different groups implies that whether an individual (city, organization) is selected for participation is decided purely by chance. Peter Rossi and Howard Freeman elaborate:

Because the resulting experimental and control groups differ from one another only by chance, whatever processes may be competing with a treatment to produce outcomes are present in the experimental and control groups to the same

extent except for chance fluctuations. For example, given randomization, persons who would be more likely to seek out the treatment if it were offered to them on a free-choice basis are equally likely to be in the experimental as in the control group. Hence, both groups have the same proportion of persons favorably predisposed to the intervention. (1985, 235)

But how is randomization accomplished? Randomization is analogous to flipping a coin, with all of those flipping heads being assigned to one group and all of those flipping tails being assigned to another. The most common ways of affecting random assignment are the following:

1. Actually flipping a coin for each subject (allowing all heads to represent one group and tails the other).
2. Throwing all the names into a hat or some other container, thoroughly mixing them, and drawing them out one at a time, allowing odd draws to represent one group and even draws the other.
3. Using a table of random numbers or

random numbers generated by a computer program.

It is important to distinguish between random assignment and random sampling. At first glance, they may appear to be identical, and in some instances they may be identical. Major differences exist, however, between the two techniques. Random sampling ensures representativeness between a sample and the population from which it is drawn. Random selection (sampling) is thus an important factor in the external validity of a study—that is, the extent to which a study’s results can be generalized beyond the sample drawn. For example, in the study presented in this chapter, Harrison McKay and colleagues evaluated a program of treatment combining nutrition, health care, and education on the cognitive ability of chronically undernourished children from around the world. Thus, the question is this (assuming that the study itself is reliable): To what degree can the impact of this program conducted in Colombia be generalized to the probable impact of similar programs on other children throughout the world?

Random assignment, as opposed to random selection, is related to the evaluation’s internal validity—that is, the extent to which the program’s impact is attributed to the treatment and no other factors.

Obviously, the best course would be to select subjects randomly and then to assign them randomly to groups once they were selected. (This discussion has digressed a bit from the discussion of the pretest-posttest control group design, but the digression is important.)

Returning to the design: Subjects are randomly assigned to an experimental group and a control group. To evaluate the effectiveness of the program, measurements are taken twice for each group. A preprogram measure is taken for each group before the introduction of the program to the experimental group. A postprogram measure is then taken after the experimental group has been exposed to (or has completed) the program.

Preprogram scores are then subtracted from postprogram scores. If the gain made by the experimental group is significantly larger than the gain made by the control group, then the researchers can conclude that the program is effective. The pretest-posttest control group design is illustrated in table 5.1. Group E is the experimental group, or the group receiving the program. Group C is the control group. The “O” indicates a test point and the “X” represents the program.

TABLE 5.1
Pretest-Posttest Control Group Design

	Pretest	Program	Posttest
Group E	O	X	O
Group C	O		O

The critical question is this: When should such a design be used? Although it is extremely powerful, this design is also costly and difficult to implement. It is not possible, for example, to withhold treatment purposely from some groups and to assign them randomly to control groups (in matters of life and death, for example). And in many cases, program participants are volunteers; there is no comparable control group (those persons who had the desire and motivation to participate in the program but who did not volunteer). Then there is the matter of cost. Experimental evaluation designs are generally more costly than other designs because of the greater amount of time required to plan and conduct the experiment and the higher level of analytical skills required for planning and undertaking the evaluation and analyzing the results.

This design is frequently used in health or employment programs. A popular, but now somewhat dated, Urban Institute publication, *Practical Program Evaluation for State and Local Governments*, documents conditions under which such experimental designs are likely to be appropriate for state and local governments (Hatry, Winnie, and Fisk 1981,

107–15). The more significant conditions include the following:

1. *There is likely to be a high degree of ambiguity as to whether outcomes were caused by the program if some other evaluation design is used.* The design is appropriate when the findings obtained through the use of a less powerful design may be criticized for not causing the results. Take a hypothetical example of a medical experiment involving a cold remedy. If no control group was used, would not critics of the experiment ask, “How do we know that the subject would not have recovered from the cold in the same amount of time without the pill?”
2. *Some citizens can be given different services from others without significant danger or harm.* The experimental design may be used if public officials and the evaluators agree that the withdrawn or nonprovision of a service or program is not likely to have harmful effects. An example might be discontinuing evening hours at a local branch library, which is not likely to harm many individuals.
3. *Some citizens can be given different services from others without violating moral and ethical standards.* Some programs, although not involving physical danger, may call for not providing services to some groups. For example, an experiment to assess the effectiveness of a counseling program for parolees might be designed in such a way that certain parolees do not receive counseling (and thus might be more likely to commit a crime and be returned to prison). This could be seen by some as unethical or immoral.
4. *There is substantial doubt about the effectiveness of a program.* If a program is not believed to be working or effective, controlled, randomized experimentation is probably the only way to settle the issue once and for all.
5. *There are insufficient resources to provide the program to all clients.* Even when a program is expected to be helpful, the resources necessary to provide it to all eligible clients may not be available. In the article in this chapter, McKay and colleagues did not have sufficient financial resources to provide the program to all chronically undernourished children in Cali, Colombia. They thus were able to evaluate the program by comparing children who had received the program with those who had not—even though it would have been desirable to provide the program to all children.
6. *The risk in funding the program without a controlled experiment is likely to be substantially greater than the cost of the experiment; the new program involves large costs and a large degree of uncertainty.* The income maintenance experiments described in chapter 2 were designed to test a new form of welfare payment. These evaluations are among the most expensive ever funded by the federal government. Yet the millions of dollars spent on the experiment are insignificant when compared to the cost of a nationwide program that did not provide the desired results.
7. *A decision to implement the program can be postponed until the experiment is completed.* Most experiments take a long time—a year or more. If there is considerable pressure (usually political) to fully implement a program, experimentation may be difficult to apply.

What all this means is that there are probably many occasions when an experimental design is not appropriate. In contrast, there are times when such a design is clearly needed. The following article, “Improving Cognitive Ability in Chronically Deprived Children,” is an example of the pretest-posttest control group design.

READING

Improving Cognitive Ability in Chronically Deprived Children

Harrison McKay • Leonardo Sinisterra • Arlene McKay
Hernando Gomez • Pascuala Lloreda

IN RECENT YEARS, social and economic planning in developing countries has included closer attention than before to the nutrition, health, and education of children of preschool age in low-income families. One basis for this, in addition to mortality and morbidity studies indicating high vulnerability at that age, (1) is information suggesting that obstacles to normal development in the first years of life, found in environments of such poverty that physical growth is retarded through malnutrition, are likely also to retard intellectual development permanently if early remedial action is not taken (2). The loss of intellectual capability, broadly defined, is viewed as especially serious because the technological character of contemporary civilization makes individual productivity and personal fulfillment increasingly contingent upon such capability. In tropical and subtropical zones of the world between 220 and 250 million children below 6 years of age live in conditions of environmental deprivation extreme enough to produce some degree of malnutrition (3); failure to act could result in irretrievable loss of future human capacity on a massive scale.

Although this argument finds widespread agreement among scientists and planners, there is uncertainty about the effectiveness of specific remedial actions. Doubts have been growing for the past decade about whether providing food, education, or health care directly to young children in poverty environments can counteract the myriad social, economic, and biological limitations to their intellectual growth. Up to 1970, when the study reported here was formulated, no definitive evidence was available to show that food

and health care provided to malnourished or "at risk" infants and young children could produce lasting increases in intellectual functioning. This was so in spite of the ample experience of medical specialists throughout the tropical world that malnourished children typically responded to nutritional recuperation by being more active physically, more able to assimilate environmental events, happier, and more verbal, all of which would be hypothesized to create a more favorable outlook for their capacity to learn (4).

In conferences and publications emphasis was increasingly placed upon the inextricable relation of malnutrition to other environmental factors inhibiting full mental development of preschool age children in poverty environments (5). It was becoming clear that, at least after the period of rapid brain growth in the first 2 years of life, when protein-calorie malnutrition could have its maximum deleterious physiological effects (6), nutritional rehabilitation and health care programs should be accompanied by some form of environmental modification of children at risk. The largest amount of available information about the potential effects of environmental modification among children from poor families pertained to the United States, where poverty was not of such severity as to make malnutrition a health issue of marked proportions. Here a large literature showed that the low intellectual performance found among disadvantaged children was environmentally based and probably was largely fixed during the preschool years (7). This information gave impetus to the belief that direct treatments, carefully designed and properly delivered to children dur-

ing early critical periods, could produce large and lasting increases in intellectual ability. As a consequence, during the 1960s a wide variety of individual, research-based preschool programs as well as a national program were developed in the United States for children from low-income families (8). Several showed positive results but in the aggregate they were not as great or as lasting as had been hoped, and there followed a widespread questioning of the effectiveness of early childhood education as a means of permanently improving intellectual ability among disadvantaged children on a large scale (9).

From pilot work leading up to the study reported here, we concluded that there was an essential issue that had not received adequate attention and the clarification of which might have tempered the pessimism: the relation of gains in intellectual ability to the intensity and duration of meliorative treatment received during different periods in the preschool years. In addition to the qualitative question of what kinds of preschool intervention, if any, are effective, attention should have been given to the question of what amount of treatment yields what amount of gain. We hypothesized that the increments in intellectual ability produced in preschool programs for disadvantaged children were subsequently lost at least in part because the programs were too brief. Although there was a consensus that longer and more intensive preschool experience could produce larger and more lasting increases, in only one study was there to be found a direct attempt to test this, and in that one sampling problems caused difficulties in interpretation (10).

As a consequence, the study reported here was designed to examine the quantitative question, with chronically undernourished children, by systematically increasing the duration of multidisciplinary treatments to levels not previously reported and evaluating results with measures directly comparable across all levels (11). This was done not only to test the hypothesis that greater amounts of

treatment could produce greater and more enduring intellectual gains but also to develop for the first time an appraisal of what results could be expected at different points along a continuum of action. This second objective, in addition to its intrinsic scientific interest, was projected to have another benefit: that of being useful in the practical application of early childhood services. Also unique in the study design was the simultaneous combination of health, nutrition, and educational components in the treatment program. With the exception of our own pilot work (12), prior studies of preschool nutritional recuperation programs had not included educational activities. Likewise, preschool education studies had not included nutritional recuperation activities, because malnutrition of the degree found in the developing countries was not characteristic of disadvantaged groups studied in the United States (13), where most of the modern early-education research had been done.

Experimental Design and Subjects

The study was carried out in Cali, Colombia, a city of nearly a million people with many problems characteristic of rapidly expanding cities in developing countries, including large numbers of families living in marginal economic conditions. Table 1 summarizes the experimental design employed. The total time available for the experiment was 3^o years, from February 1971 to August 1974. This was divided into four treatment periods of 9 months each plus interperiod recesses. Our decision to begin the study with children as close as possible to 3 years of age was based upon the 2 years of pilot studies in which treatment and measurement systems were developed for children starting at that age (14). The projected 180 to 200 days of possible attendance at treatment made each projected period similar in length to a school

TABLE 1
Basic selection and treatment variables of the groups of children in the study.

Group	N		Characteristic
	In 1971	In 1975	
T1(a)	57	49	Low SES, subnormal weight and height. One treatment period, between November 1973 and August 1974 (75 to 84 months of age) ^a
T1(b)	56	47	Low SES, subnormal weight and height. One treatment period, between November 1973 and August 1974 (75 to 84 months of age), with prior nutritional supplementation and health care
T2	64	51	Low SES, subnormal weight and height. Two treatment periods, between November 1972 and August 1974 (63 to 84 months of age)
T3	62	50	Low SES, subnormal weight and height. Three treatment periods, between December 1971 and August 1974 (52 to 84 months of age)
T4	62	51	Low SES, subnormal weight and height. Four treatment periods, between February 1971 and August 1974 (42 to 84 months of age)
HS	38	30	High SES. Untreated, but measured at the same points as groups T1–T4
T0	116	72	Low SES, normal weight and height. Untreated

a. SES is family socioeconomic status.

year in Colombia, and the end of the fourth period was scheduled to coincide with the beginning of the year in which the children were of eligible age to enter first grade.

With the object of having 60 children initially available for each treatment group (in case many should be lost to the study during the 3° year period), approximately 7500 families living in two of the city's lowest-income areas were visited to locate and identify all children with birth dates between 1 June and 30 November 1967, birth dates that would satisfy primary school entry requirements in 1974. In a second visit to the 733 families with such children, invitations were extended to have the children medically examined. The families of 518 accepted, and each child received a clinical examination, anthropometric measurement, and screening for serious neurological dysfunctions. During a third visit to these families, interviews and observations were conducted to determine living conditions, economic resources, and age, education, and

occupations of family members. At this stage the number of potential subjects was reduced by 69 (to 449), because of errors in birth date, serious neurological or sensory dysfunctions, refusal to participate further, or removal from the area.

Because the subject loss due to emigration during the 4 months of preliminary data gathering was substantial, 333 children were selected to assure the participation of 300 at the beginning of treatment; 301 were still available at that time, 53 percent of them male. Children selected for the experiment from among the 449 candidates were those having, first, the lowest height and weight for age; second, the highest number of clinical signs of malnutrition (15); and third, the lowest per capita family income. The second and third criteria were employed only in those regions of the frequency distributions where differences among the children in height and weight for age were judged by the medical staff to lack biological significance. Figure 1 shows these frequency distributions

and includes scales corresponding to percentiles in a normal population (16).

The 116 children not selected were left untreated and were not measured again until 4 years later, at which point the 72 still living in the area and willing once again to collaborate were reincorporated into the longitudinal study and measured on physical growth and cognitive development at the same time as the selected children, beginning at 7 years of age. At 3 years of age these children did not show abnormally low weight for age or weight for height.

In order to have available a set of local reference standards for “normal” physical and psychological development, and not depend solely upon foreign standards, a group of children (group HS) from families with high socioeconomic status, living in the same city and having the same range of birth dates as the experimental group, was included in the study. Our assumption was that, in regard to available economic resources, housing, food, health care, and educational opportunities, these children had the highest probability of full intellectual and physical development of any group in the society. In relation to the research program they remained untreated, receiving only medical and psychological assessment at the same intervals as the treated children, but the majority were attending the best private preschools during the study. Eventually 63 children were recruited for group HS, but only the 38 noted in Table 1 were available at the first psychological testing session in 1971.

Nearly all the 333 children selected for treatment lived in homes distributed throughout an area of approximately 2 square kilometers. This area was subdivided into 20 sectors in such a way that between 13 and 19 children were included in each sector. The sectors were ranked in order of a standardized combination of average height and weight for age and per capita family income of the children. Each of

the first five sectors in the ranking was assigned randomly to one of five groups. This procedure was followed for the next three sets of five sectors, yielding four sectors for each group, one from each of four strata. At this point the groups remained unnamed; only as each new treatment period was to begin was a group assigned to it and families in the sectors chosen so informed. The children were assigned by sectors instead of individually in order to minimize social interaction between families in different treatment groups and to make daily transportation more efficient (17). Because this “lottery” system was geographically based, all selected children who were living in a sector immediately prior to its assignment were included in the assignment and remained in the same treatment group regardless of further moves. In view of this process, it must be noted that the 1971 N’s reported for the treatment groups in Table 1 are retrospective figures, based upon a count of children then living in sectors assigned later to treatment groups. Table 1 also shows the subject loss, by group, between 1971 and 1975. The loss of 53 children—18 percent—from the treatment groups over 4 years was considerably less than expected. Two of these children died and 51 emigrated from Cali with their families; on selection variables they did not differ to a statistically significant degree from the 248 remaining.

A longitudinal study was begun, then, with groups representing extreme points on continua of many factors related to intellectual development, and with an experimental plan to measure the degree to which children at the lower extreme could be moved closer to those of the upper extreme as a result of combined treatments of varying durations. Table 2 compares selected (T1-T4), not selected (T0), and reference (HS) groups on some of the related factors, including those used for selecting children for participation in treatment.

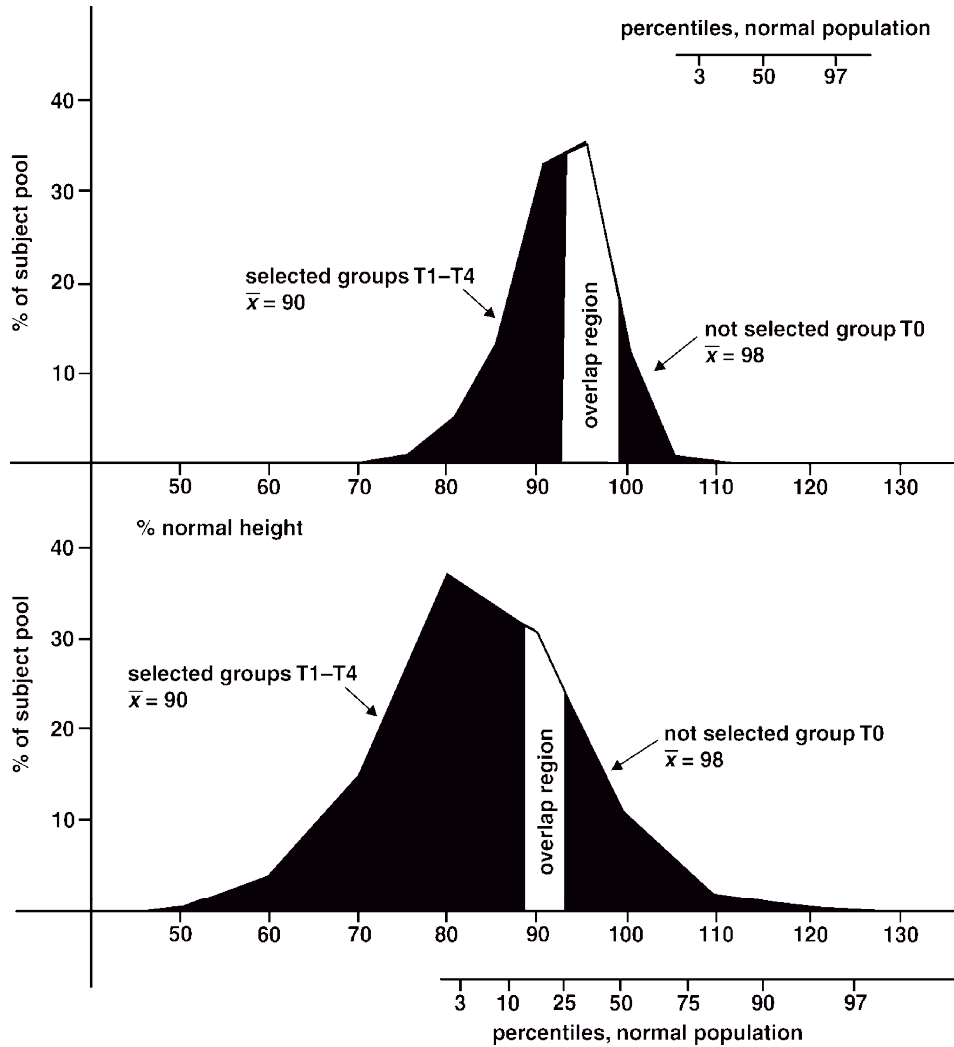


FIGURE 1

Frequency distributions of height and weight (as percent of normal for age) of the subject pool of 449 children available in 1970, from among whom 333 were selected for treatment groups. A combination of height and weight was the first criterion; the second and third criteria, applied to children in the overlap regions, were clinical signs of malnutrition and family income. Two classification systems for childhood malnutrition yield the following description of the selected children: 90 percent nutritionally “stunted” at 3 years of age and 35 percent with evidence of “wasting”; 26 percent with “second degree” malnutrition, 54 percent with “first degree,” and 16 percent “low normal.” (16)

Treatments

The total number of treatment days per period varied as follows: period 1, 180 days; period 2, 185; period 3, 190; period 4, 172. A fire early in period 4 reduced the time available owing to the necessity of terminating the study before the opening of primary school. The original objective was to have each succeeding period at least as long as the preced-

ing one in order to avoid reduction in intensity of treatment. The programs occupied 6 hours a day 5 days a week, and attendance was above 95 percent for all groups; hence there were approximately 1040, 1060, 1080, and 990 hours of treatment per child per period from period 1 to period 4, respectively. The total number of hours of treatment per group, then, were as follows: T4, 4170 hours;

T3, 3130 hours; T2, 2070 hours; T1 (a and b), 990 hours.

In as many respects as possible, treatments were made equivalent between groups within each period. New people, selected and trained as child-care workers to accommodate the periodic increases in numbers of children, were combined with existing personnel and distributed in such a way that experience, skill, and familiarity with children already treated were equalized for all groups, as was the adult-child ratio. Similarly, as new program sites were added, children rotated among them so that all groups occupied all sites equal lengths of time. Except for special care given to the health and nutritional adaptation of each newly entering group during the initial weeks, the same systems in these treatments were applied to all children within periods.

An average treatment day consisted of 6 hours of integrated health, nutritional, and educational activities, in which approximately 4 hours were devoted to education and 2 hours to health, nutrition, and hygiene. In practice, the nutrition and health care provided opportunities to reinforce many aspects of the educational curriculum, and time in the education program was used to reinforce recommended hygienic and food consumption practices.

The nutritional supplementation program was designed to provide a minimum of 75 percent of recommended daily protein and calorie allowances, by means of low-cost foods available commercially, supplemented with vitamins and minerals, and offered ad libitum three times a day. In the vitamin and mineral supplementation, special attention was given to vitamin A, thiamin, riboflavin, niacin, and iron, of which at least 100 percent of recommended dietary allowance was provided (18).

The health care program included daily observation of all children attending the treatment center, with immediate pediatric attention to those with symptoms reported by the parents or noted by the health and

education personnel. Children suspected of an infectious condition were not brought into contact with their classmates until the danger of contagion had passed. Severe health problems occurring during weekends or holidays were attended on an emergency basis in the local university hospital.

The educational treatment was designed to develop cognitive processes and language, social abilities, and psychomotor skills, by means of an integrated curriculum model. It was a combination of elements developed in pilot studies and adapted from other programs known to have demonstrated positive effects upon cognitive development (19). Adapting to developmental changes in the children, its form progressed from a structured day divided among six to eight different directed activities, to one with more time available for individual projects. This latter form, while including activities planned to introduce new concepts, stimulate verbal expression, and develop motor skills, stressed increasing experimentation and decision taking by the children. As with the nutrition and health treatments during the first weeks of each new period, the newly entering children received special care in order to facilitate their adaptation and to teach the basic skills necessary for them to participate in the program. Each new period was conceptually more complex than the preceding one, the last ones incorporating more formal reading, writing, and number work.

Measures of Cognitive Development

There were five measurement points in the course of the study: (i) at the beginning of the first treatment period; (ii) at the end of the first treatment period; (iii) after the end of the second period, carrying over into the beginning of the third; (iv) after the end of the third period, extending into the fourth; and (v) following the fourth treatment pe-

TABLE 2

Selection variables and family characteristics of study groups in 1970 (means). All differences between group HS and groups T1–T4 are statistically significant ($P < .01$) except age of parents. There are no statistically significant differences among groups T1–T4. There are statistically significant differences between group T0 and combined groups T1–T4 in height and weight (as percent of normal), per capita income and food expenditure, number of family members and children, and rooms per child; and between group T0 and group HS on all variables except age of parents and weight.

Variable	Group		
	T1–T4	T0	HS
Height as percent of normal for age	90	98	101
Weight as percent of normal for age	79	98	102
Per capita family income as percent of group HS	5	7	100
Per capita food expenditure in family as percent of group HS	15	22	100
Number of family members	7.4	6.4	4.7
Number of family under 15 years of age	4.8	3.8	2.4
Number of play/sleep rooms per child	0.3	0.5	1.6
Age of father	37	37	37
Age of mother	31	32	31
Years of schooling, father	3.6	3.7	14.5
Years of schooling, mother	3.5	3.3	10.0

riod. For the purpose of measuring the impact of treatment upon separate components of cognitive development, several short tests were employed at each measurement point, rather than a single intelligence test. The tests varied from point to point, as those only applicable at younger ages were replaced by others that could be continued into primary school years. At all points the plan was to have tests that theoretically measured adequacy of language usage, immediate memory, manual dexterity and motor control, information and vocabulary, quantitative concepts, spatial relations, and logical thinking, with a balance between verbal and nonverbal production. Table 3 is a list of tests applied at each measurement point. More were applied than are listed; only those employed at two or more measurement points and having items that fulfilled the criteria for the analysis described below are included.

Testing was done by laypersons trained and supervised by professional psychologists. Each new test underwent a 4 to 8 month developmental sequence which included an

initial practice phase to familiarize the examiners with the format of the test and possible difficulties in application. Thereafter, a series of pilot studies were conducted to permit the modification of items in order to attain acceptable levels of difficulty, reliability, and ease of application. Before each measurement point, all tests were applied to children not in the study until adequate inter-tester reliability and standardization of application were obtained. After definitive application at each measurement point, all tests were repeated on a 10 percent sample to evaluate test-retest reliability. To protect against examiner biases, the children were assigned to examiners randomly and no information was provided regarding treatment group or nutritional or socioeconomic level. (The identification of group HS children was, however, unavoidable even in the earliest years, not only because of their dress and speech but also because of the differences in their interpersonal behavior.) Finally, in order to prevent children from being trained specifically to perform well on test items, the two functions of intervention and evaluation were separated as far as possible. We intention-

ally avoided, in the education programs, the use of materials or objects from the psychological tests. Also, the intervention personnel had no knowledge of test content or format, and neither they nor the testing personnel were provided with information about group performance at any of the measurement points.

Data Analysis

The data matrix of cognitive measures generated during the 44-month interval between the first and last measurement points entailed evaluation across several occasions by means of a multivariate vector of observations. A major problem in the evaluation procedure, as seen in Table 3, is that the tests of cognitive development were not the same at every measurement point. Thus the re-

sponse vector was not the same along the time dimension. Initially, a principal component approach was used, with factor scores representing the latent variables (20). Although this was eventually discarded because there was no guarantee of factor invariance across occasions, the results were very similar to those yielded by the analyses finally adopted for this article. An important consequence of these analyses was the finding that nearly all of the variation could be explained by the first component (21), and under the assumption of unidimensionality cognitive test items were pooled and calibrated according to the psychometric model proposed by Rasch (22) and implemented computationally by Wright (23). The technique employed to obtain the ability estimates in Table 4 guarantees that the same latent trait is being reflected in these estimates (24). Consequently, the growth curves in Fig. 2 are interpreted as representing “general cognitive ability” (25).

Table 5 shows correlations between pairs of measurement points of the ability estimates of all children included in the two points. The correspondence is substantial, and the matrix exhibits the “simplex” pattern expected in psychometric data of this sort (26). As the correlations are not homogeneous, a test for diagonality in the transformed error covariance matrix was carried out, and the resulting chi-square value led to rejection of a mixed model assumption. In view of this, Bock’s multivariate procedure (27), which does not require constant correlations, was employed to analyze the differences among groups across measurement points. The results showed a significant groups-by-occasions effect, permitting rejection of the hypothesis of parallel profiles among groups. A single degree-of-freedom decomposition of this effect showed that there were significant differences in every possible Helmert contrast. Stepdown tests indicated that all components were required in describing profile differences.

TABLE 3

Tests of cognitive ability applied at different measurement points (see text) between 43 and 87 months of age. Only tests that were applied at two adjacent points and that provided items for the analysis in Table 4 are included. The unreferenced tests were constructed locally.

Test	Measurement points
Understanding complex commands	1,2
Figure tracing	1, 2, 3
Picture vocabulary	1, 2, 3
Intersensory perception (33)	1, 2, 3
Colors, numbers, letters	1, 2, 3
Use of prepositions	1, 2, 3
Block construction	1, 2, 3
Cognitive maturity (34)	1, 2, 3, 4
Sentence completion (35)	1, 2, 3, 4
Memory for sentences (34)	1, 2, 3, 4, 5
Knox cubes (36)	1, 2, 3, 4, 5
Geometric drawings (37)	3, 4
Arithmetic (38, 39)	3, 4, 5
Mazes (40)	3, 4, 5
Information (41)	3, 4, 5
Vocabulary (39)	3, 4, 5
Block design (42)	4, 5
Digit memory (43)	4, 5
Analogies and similarities (44)	4, 5
Matrices (45)	4, 5
Visual classification	4, 5

TABLE 4

Scaled scores on general cognitive ability, means and estimated standard errors, of the four treatment groups and group HS at five testing points.

Group	N	Average age at testing (months)				
		43	49	63	77	87
<i>Mean score</i>						
HS	28	-0.11	.39	2.28	4.27	4.89
T4	50	-1.82 ^a	.21	1.80	3.35	3.66
T3	47	-1.72	-1.06	1.64	3.06	3.35
T2	49	-1.94	-1.22	.30 ^b	2.61	3.15
T1	90	-1.83	-1.11	.33	2.07	2.73
<i>Estimated standard error</i>						
HS	28	.192	.196	.166	.191	.198
T4	50	.225	.148	.138	.164	.152
T3	47	.161	.136	.103	.123	.120
T2	49	.131	.132	.115	.133	.125
T1	90	.110	.097	.098	.124	.108
<i>Standard deviation</i>						
All groups		1.161	1.153	1.169	1.263	1.164

a. Calculated from 42 percent sample tested prior to beginning of treatment.

b. Calculated from 50 percent sample tested prior to beginning of treatment.

The data in Table 4, plotted in Fig. 2 with the addition of dates and duration of treatment periods, are based upon the same children at all measurement points. These are children having complete medical, socioeconomic, and psychological test records. The discrepancies between the 1975 *N*'s in Table 1 and the *N*'s in Table 4 are due to the fact that 14 children who were still participating in the study in 1975 were excluded from the analysis because at least one piece of information was missing, a move made to facilitate correlational analyses. Between 2 percent (T4) and 7 percent (HS) were excluded for this reason.

TABLE 5

Correlation of ability scores across measurement points

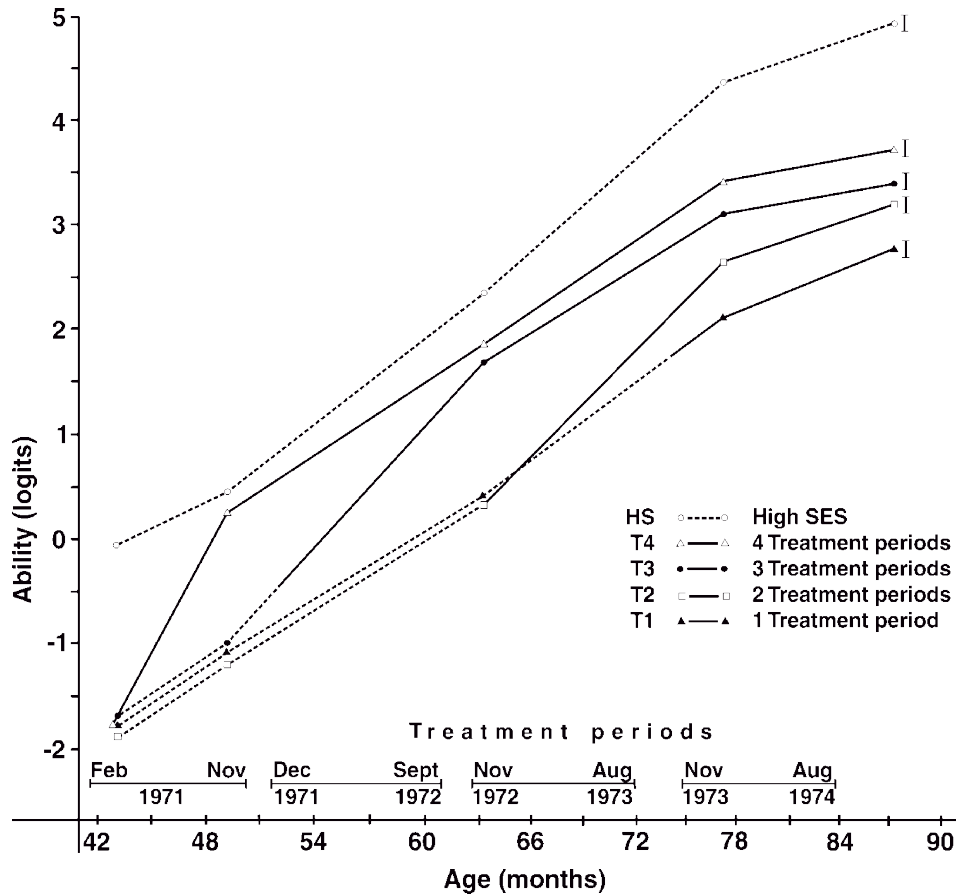
Measurement points	1	2	3	4	5
1	—	.78	.68	.54	.48
2		—	.80	.66	.59
3			—	.71	.69
4				—	.76
5					—

For all analyses, groups T1(a) and T1(b) were combined into group T1 because the prior nutritional supplementation and health care provided group T1(b) had not been found to produce any difference between the two groups. Finally, analysis by sex is not included because a statistically significant difference was found at only one of the five measurement points.

Relation of Gains to Treatment

The most important data in Table 4 and Fig. 2 are those pertaining to cognitive ability scores at the fifth testing point. The upward progression of mean scores from T1 to T4 and the nonoverlapping standard errors, except between T2 and T3, generally confirm that the sooner the treatment was begun the higher the level of general cognitive ability reached by age 87 months. Another interpretation of the data could be that the age at which treatment began was a determining factor independent of amount of time in treatment.

It can be argued that the level of cognitive development which the children reached at 7

**FIGURE 2**

Growth of general cognitive ability of the children from age 43 months to 87 months, the age at the beginning of primary school. Ability scores are scaled sums of test items correct among items common to proximate testing points. The solid lines represent periods of participation in a treatment sequence, and brackets to the right of the curves indicate ± 1 standard error of the corresponding group means at the fifth measurement point. At the fourth measurement point there are no overlapping standard errors; at earlier measurement points there is overlap only among obviously adjacent groups (see Table 4). Group T0 was tested at the fifth measurement point but is not represented in this figure, or in Table 2, because its observed low level of performance could have been attributed to the fact that this was the first testing experience of the group T0 children since the neurological screening 4 years earlier.

years of age depended upon the magnitude of gains achieved during the first treatment period in which they participated, perhaps within the first 6 months, although the confounding of age and treatment duration in the experimental design prohibits conclusive testing of the hypothesis. The data supporting this are in the declining magnitude of gain, during the first period of treatment attended, at progressively higher ages of entry into the program. Using group T1 as an untreated baseline until it first entered treatment, and

calculating the difference in gains (28) between it and groups T4, T3, and T2 during their respective first periods of treatment, we obtain the following values: group T4, 1.31; group T3, 1.26; and group T2, .57. When calculated as gains per month between testing periods, the data are the following: T4, .22; T3, .09; and T2, .04. This suggests an exponential relationship. Although, because of unequal intervals between testing points and the overlapping of testing durations with treatment periods, this latter relationship must be

viewed with caution, it is clear that the older the children were upon entry into the treatment programs the less was their gain in cognitive development in the first 9 months of participation relative to an untreated baseline.

The lack of a randomly assigned, untreated control group prevents similar quantification of the response of group T1 to its one treatment period. If group HS is taken as the baseline, the observed gain of T1 is very small. The proportion of the gap between group HS and group T1 that was closed during the fourth treatment period was 2 percent, whereas in the initial treatment period of each of the other groups the percentages were group T4, 89; group T3, 55; and group T2, 16. That the progressively declining responsiveness at later ages extends to group T1 can be seen additionally in the percentages of gap closed between group T4 and the other groups during the first treatment period of each of the latter: group T3, 87; group T2, 51; and group T1, 27.

Durability of Gains

Analysis of items common to testing points five and beyond has yet to be done, but the data contained in Fig. 3, Stanford-Binet intelligence quotients at 8 years of age, show that the relative positions of the groups at age 7 appear to have been maintained to the end of the first year of primary school. Although the treated groups all differ from each other in the expected direction, generally the differences are not statistically significant unless one group has had two treatment periods more than another. A surprising result of the Stanford-Binet testing was that group T0 children, the seemingly more favored among the low-income community (see Table 2), showed such low intelligence quotients; the highest score in group T0 (IQ = 100) was below the mean of group HS, and the lowest group HS score (IQ = 84) was above the mean of group T0. This further confirms that

the obstacles to normal intellectual growth found in conditions of poverty in which live large segments of the population are very strong. It is possible that this result is due partly to differential testing histories, despite the fact that group T0 had participated in the full testing program at the preceding fifth measurement point, and that this was the first Stanford-Binet testing for the entire group of subject children.

The difference between groups T0 and T1 is in the direction of superiority of group T1 ($t = 1.507$, $P < .10$). What the IQ of group T1 would have been without its one treatment period is not possible to determine except indirectly through regression analyses with other variables, but we would expect it to have been lower than T0's, because T0 was significantly above T1 on socioeconomic and anthropometric correlates of IQ (Table 2). Also, T1 was approximately .30 standard deviation below T0 at 38 months of age on a cognitive development factor of a preliminary neurological screening test applied in 1970, prior to selection. Given these data and the fact that at 96 months of age there is a difference favoring group T1 that approaches statistical significance, we conclude not only that group T1 children increased in cognitive ability as a result of their one treatment period (although very little compared to the other groups) but also that they retained the increase through the first year of primary school.

An interesting and potentially important characteristic of the curves in Fig. 3 is the apparent increasing bimodality of the distribution of the groups with increasing length of treatment, in addition to higher means and upward movement of both extremes. The relatively small sample sizes and the fact that these results were found only once make it hazardous to look upon them as definitive. However, the progression across groups is quite uniform and suggests that the issue of individual differential response to equivalent treatment should be studied more carefully.

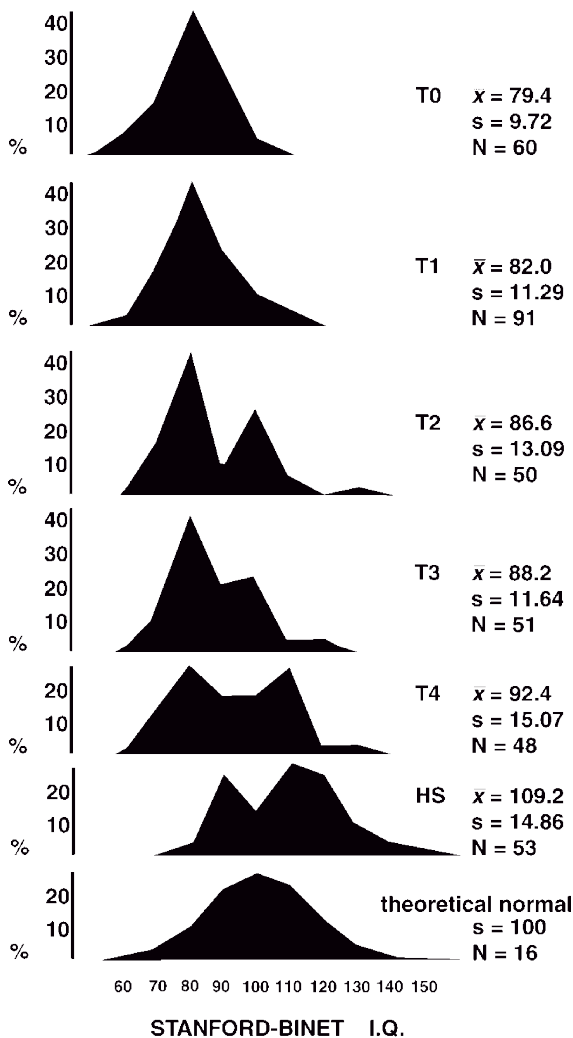
Social Significance of Gains

Group HS was included in the study for the purpose of establishing a baseline indicating what could be expected of children when conditions for growth and development were

FIGURE 3

Mean scores on the Stanford-Binet Intelligence Test at 8 years of age. Groups T0–T4 had had one year of primary school. Group HS children had attended pre-school and primary schools for up to five consecutive years prior to this testing point. Mental age minus chronological age is as follows:

Group T0	- 18 months
Group T1	- 15 months
Group T2	- 11 months
Group T3	- 9 months
Group T4	- 5 months
Group HS	+ 10 months



optimal. In this way the effectiveness of the treatment could be evaluated from a frame of reference of the social ideal. It can be seen in Table 4 that group HS increased in cognitive ability at a rate greater than the baseline group T1 during the 34 months before T1 entered treatment. This is equivalent to, and confirms, the previously reported relative decline in intelligence among disadvantaged children (29, p. 258). Between the ages of 4 and 6 years, group HS children passed through a period of accelerated development that greatly increased the distance between them and all the treatment groups. The result, at age 77 months, was that group T4 arrived at a point approximately 58 percent of the distance between group HS and the untreated baseline group T1, group T3 arrived at 45 percent, and group T2 at 24 percent. Between 77 and 87 months, however, these differences appear to have diminished, even taking into account that group T1 entered treatment during this period. In order for these percentages to have been maintained, the baseline would have had to remain essentially unchanged. With respect to overall gains from 43 months to 87 months, the data show that reduction of the 1.5 standard deviation gap found at 43 months of age between group HS and the treated children required a duration and intensity of treatment at least equal to that of group T2; the group HS overall growth of 5.00 units of ability is less than that of all groups except T1.

As noted, group HS was not representative of the general population, but was a sample of children intentionally chosen from a subgroup above average in the society in characteristics favorable to general cognitive development. For the population under study, normative data do not exist; the “theoretical normal” distribution shown in Fig. 3 represents the U.S. standardization group of 1937 (30). As a consequence, the degree to which the treatments were effective in closing the gap between the disadvantaged children and what could be de-

scribed as an acceptable level cannot be judged. It is conceivable that group HS children were developing at a rate superior to that of this hypothetical normal. If that was the case, the gains of the treated children could be viewed even more positively.

Recent studies of preschool programs have raised the question whether differences between standard intellectual performance and that encountered in disadvantaged children represent real deficits or whether they reflect cultural or ethnic uniquenesses. This is a particularly relevant issue where disadvantaged groups are ethnically and linguistically distinct from the dominant culture (29, pp. 262–72; 31). The historical evolution of differences in intellectual ability found between groups throughout the world is doubtless multidimensional, with circumstances unique to each society or region, in which have entered religious, biological, and other factors in different epochs, and thus the simple dichotomy of culture uniqueness versus deprivation is only a first approximation to a sorely needed, thorough analysis of antecedents and correlates of the variations. Within the limits of the dichotomy, however, the evidence with regard to the children in our study suggests that the large differences in cognitive ability found between the reference group and the treated groups in 1971 should be considered as reflecting deficits rather than divergent ethnic identities. Spanish was the language spoken in all the homes, with the addition of a second language in some group HS families. All the children were born in the same city sharing the same communication media and popular culture and for the most part the same religion. Additionally, on tests designed to maximize the performance of the children from low-income families by the use of objects, words, and events typical in their neighborhoods (for example, a horse-drawn cart in the picture vocabulary test), the difference between them and group HS was still approximately 1.50 standard deviations at 43 months

of age. Thus it is possible to conclude that the treated children's increases in cognitive ability are relevant to them in their immediate community as well as to the ideal represented by the high-status reference group. This will be more precisely assessed in future analyses of the relation of cognitive gains to achievement in primary school.

Conclusions

The results leave little doubt that environmental deprivation of a degree severe enough to produce chronic undernutrition manifested primarily by stunting strongly retards general cognitive development, and that the retardation is less amenable to modification with increasing age. The study shows that combined nutritional, health, and educational treatments between 3° and 7 years of age can prevent large losses of potential cognitive ability, with significantly greater effect the earlier the treatments begin. As little as 9 months of treatment prior to primary school entry appears to produce significant increases in ability, although small compared to the gains of children receiving treatment lasting two, three, and four times as long. Continued study will be necessary to ascertain the long-range durability of the treatment effects, but the present data show that they persist at 8 years of age.

The increases in general cognitive ability produced by the multiform preschool interventions are socially significant in that they reduce the large intelligence gap between children from severely deprived environments and those from favored environments, although the extent to which any given amount of intervention might be beneficial to wider societal development is uncertain (32). Extrapolated to the large number of children throughout the world who spend their first years in poverty and hunger, however, even the smallest increment resulting from one 9-month treatment period could constitute an

important improvement in the pool of human capabilities available to a given society.

References and Notes

1. D. B. Jelliffe, *Infant Nutrition in the Tropics and Subtropics* (World Health Organization, Geneva, 1968); C. D. Williams, in *Preschool Child Malnutrition* (National Academy of Sciences, Washington, D. C., 1966), pp. 3–8; N. S. Scrimshaw, in *ibid.*, pp. 63–73.
2. N. S. Scrimshaw and J. E. Gordon, Eds., *Malnutrition, Learning and Behavior* (MIT Press, Boston, 1968), especially the articles by J. Cravioto and E. R. DeLicardie, F. Monckeberg, and M. B. Stoch and P. M. Smythe; J. Cravioto, in *Preschool Child Malnutrition* (National Academy of Sciences, Washington, D.C., 1966), pp. 74–84; M. Winick and P. Rosso, *Pediatr. Res.* 3, 181 (1969); L. M. Brockman and H. N. Ricciuti, *Dev. Psychol.* 4, 312 (1971). As significant as the human studies was the animal research of R. Barnes, J. Cowley, and S. Franková, all of whom summarize their work in *Malnutrition, Learning and Behavior*.
3. A. Berg, *The Nutrition Factor* (Brookings Institution, Washington, D.C., 1973), p. 5.
4. In addition to the authors' own experience, that of pediatricians and nutrition specialists in Latin America, Africa, and Asia is highly uniform in this respect. In fact, a generally accepted rule in medical care of malnourished children is that improvement in psychological response is the first sign of recovery.
5. D. Kallen, Ed., *Nutrition, Development and Social Behavior*. DHEW Publication No. (NIH) 73–242 (National Institutes of Health, Washington, D.C., 1973); S. L. Manocha, *Malnutrition and Retarded Human Development* (Thomas, Springfield, Ill., 1972). pp. 132–165.
6. J. Dobbing, in *Malnutrition, Learning and Behavior*, N. S. Scrimshaw and J. E. Gordon, Eds. (MIT Press, Boston, 1968), p. 181.
7. B. S. Bloom, *Stability and Change in Human Characteristics* (Wiley, New York, 1964); J. McV. Hunt, *Intelligence and Experience* (Ronald, New York, 1961); M. Deutsch et al., *The Disadvantaged Child* (Basic Books, New York, 1967). As with nutrition studies, there was also a large amount of animal research on early experience, such as that of J. Denenberg and G. Karas [*Science* 130, 629 (1959)] showing the vulnerability of the young organism and that early environmental effects could last into adulthood.
8. M. Pines, *Revolution in Learning* (Harper & Row, New York, 1966); R. O. Hess and R. M. Bear, Eds., *Early Education* (Aldine, Chicago, 1968).
9. The 1969–1970 debate over the effectiveness of early education can be found in M. S. Smith and J. S. Bissel, *Harv. Educ. Rev.* 40, 51 (1970); V. G. Cicirelli, J. W. Evans, J. S. Schiller, *ibid.* p. 105; D. T. Campbell and A. Erlebacher, in *Disadvantaged Child*, J. Hellmuth, Ed. (Brunner/Mazel, New York, 1970), vol. 3; *Environment, Heredity and Intelligence* (Harvard Educational Review, Cambridge, Mass., 1969); F. D. Horowitz and L. Y. Paden, in *Review of Child Development Research*, vol. 3, *Child Development and Social Policy*, B. M. Caldwell and H. N. Ricciuti, Eds. (Univ. of Chicago Press, Chicago, 1973), pp. 331–402; T. Kellaghan, *The Evaluation of a Preschool Programme for Disadvantaged Children* (Educational Research Centre, St. Patrick's College, Dublin, 1975), pp. 21–33; U. Bronfenbrenner, *Is Early Intervention Effective?*, DHEW Publication No. (OHD) 74–25 (Office of Human Development, Department of Health, Education, and Welfare, Washington, D.C., 1974).
10. S. Gray and R. Klaus, in *Intelligence: Some Recurring Issues*, L. E. Tyler, Ed. (Van Nostrand Reinhold, New York, 1969), pp. 187–206.
11. This study sprang from prior work attempting to clarify the relationship of moderate malnutrition to mental retardation, in which it became evident that malnutrition of first and second degree might be only one of many correlated environmental factors found in poverty contributing to deficit in cognitive performance. We selected children for treatment with undernutrition as the first criterion not to imply that this was the major critical factor in cognitive development, but to be assured that we were dealing with children found typically in the extreme of poverty in the developing world, permitting generalization of our results to a population that is reasonably well defined in pediatric practice universally. Figure 1 allows scientists anywhere to directly compare their population with ours on criteria that, in our experience, more reliably reflect the sum total of chronic environment deprivation than any other measure available.
12. H. McKay, A. McKay, L. Sinisterra, in *Nutrition, Development and Social Behavior*, D. Kallen, Ed., DHEW Publ. No. (NIH) 73–242 (National Institutes of Health, Washington, D.C., 1973); S. Franková, in *Malnutrition, Learning and Behavior*, N.S. Scrimshaw and J. E. Gordon, Eds. (MIT Press, Cambridge, Mass., 1968). Franková, in addition to showing in her studies with rats that malnutrition caused behavioral abnormalities, also was the first to suggest that the effects of malnutrition could be modified through stimulation.
13. *First Health and Nutrition Examination Survey, United States, 1971–72*, DHEW Publication No. (HRA) 75–1223 (Health Resources Administra-

- tion, Department of Health, Education, and Welfare, Rockville, Md., 1975).
14. H. McKay, A. McKay, L. Sinisterra, *Stimulation of Cognitive and Social Competence in Preschool Age Children Affected by the Multiple Deprivations of Depressed Urban Environments* (Human Ecology Research Foundation, Cali, Colombia, 1970).
 15. D. B. Jelliffe, *The Assessment of the Nutritional Status of the Community* (World Health Organization, Geneva, 1966), pp. 10–96.
 16. In programs that assess community nutritional conditions around the world, standards widely used for formal growth of preschool age children have been those from the Harvard School of Public Health found in *Textbook of Pediatrics*, W. E. Nelson, V. C. Vaughan, R. J. McKay, Eds. (Saunders, Philadelphia, ed. 9, 1969), pp. 15–57. That these were appropriate for use with the children studied here is confirmed by data showing that the “normal” comparison children (group HS) were slightly above the median values of the standards at 3 years of age. Height for age, weight for age, and weight for height of all the study children were compared with this set of standards. The results, in turn, were the basis for the Fig. 1 data, including the “stunting” (indication of chronic, or past, undernutrition) and “wasting” (indication of acute, or present, malnutrition) classifications of J. C. Waterlow and R. Rutishauser [in *Early Malnutrition and Mental Development*, J. Cravioto, L. Hambreaus, B. Vahlquist, Eds. (Almqvist & Wiksell, Uppsala, Sweden, 1974), pp. 13–26]. The classification “stunted” included the children found in the range from 75 to 95 percent of height for age. The “wasting” classification included children falling between 75 and 90 percent of weight for height. The cutoff points of 95 percent height for age and 90 percent weight for height were, respectively, 1.8 and 1.2 standard deviations (S.D.) below the means of group HS, while the selected children’s means were 3 S.D. and 1 S.D. below group HS. The degree of malnutrition, in accordance with F. Gomez, R. Ramos-Galvan, S. Frenk, J. M. Cravioto, J. M. Chavez, and J. Vasquez [J. Trop. Pediatr. 2, 77 (1956)] was calculated with weight for age less than 75 percent as “second degree” and 75 to 85 percent as “first degree.” We are calling low normal a weight for age between 85 and 90 percent. The 75, 85, and 90 percent cutoff points were 2.6, 1.7, and 1.3 S.D. below the mean of group HS, respectively, while the selected children’s mean was 2.3 S.D. below that of group HS children. Examination of combinations of both height and weight to assess nutritional status follows the recommendations of the World Health Organization expert committee on nutrition, found in *FAO/WHO Technical Report Series No. 477* (World Health Organization, Geneva, 1972), pp. 36–48 (Spanish language edition). In summary, from the point of view of both the mean values and the severity of deficit in the lower extreme of the distributions, it appears that the group of selected children, at 3 years of age, can be characterized as having had a history of chronic undernutrition rather than suffering from acute malnutrition at the time of initial examination.
 17. The data analyses in this article use individuals as the randomization unit rather than sectors. To justify this, in view of the fact that random assignment of children was actually done by sectors, a nested analysis of variance was performed on psychological data at the last data point, at age 7, to examine the difference between mean-square for variation (MS) between sectors within treatment and MS between subjects within treatments within sectors. The resulting insignificant F-statistic ($F = 1.432$, d.f. 15,216) permits such analyses.
 18. Food and Nutrition Board, National Research Council, *Recommended Dietary Allowances*, (National Academy of Sciences, Washington, D.C., 1968).
 19. D. B. Weikart acted as a principal consultant on several aspects of the education program; D. B. Weikart, L. Rogers, C. Adcock, D. McClelland [*The Cognitively Oriented Curriculum* (ERIC/National Association for the Education of Young Children, Washington, 1971)] provided some of the conceptual framework. The content of the educational curriculum included elements described in the works of C. S. Lavatelli, *Piaget’s Theory Applied to an Early Childhood Curriculum* (Center for Media Development, Boston, 1970); S. Smilansky, *The Effects of Sociodramatic Play on Disadvantaged Preschool Children* (Wiley, New York, 1968); C. Bereiter and S. Englemann, *Teaching Disadvantaged Children in the Preschool* (Prentice-Hall, Englewood Cliffs, N.J., 1969); R. G. Stauffer, *The Language-Experience Approach to the Teaching of Reading* (Harper & Row, New York, 1970); R. Van Allen and C. Allen, *Language Experiences in Early Childhood* (Encyclopaedia Britannica Press, Chicago, 1969); S. Ashton-Warner, *Teacher* (Bantam, New York, 1963); R. C. Orem, *Montessori for the Disadvantaged Children in the Preschool* (Capricorn, New York, 1968); and M. Montessori, *The Discovery of the Child* (Ballantine, New York, 1967).
 20. M. McKay, L. Sinisterra, A. McKay, H. Gomez, P. Lloreda, J. Korgi, A. Dow, in *Proceedings of the Tenth International Congress of Nutrition* (International Congress of Nutrition, Kyoto, 1975), chap. 7.
 21. Although the “Scree” test of R. B. Cattell [*Multivar. Behav. Res.* 2, 245 (1966)] conducted

on the factor analyses at each measurement point clearly indicated a single factor model, there does exist the possibility of a change in factorial content that might have affected the differences between group HS children and the others at later measurement periods. New analysis procedures based upon a linear structural relationship such as suggested by K. G. Jöreskog and D. Sörbom [in *Research Report 76-1* (Univ. of Uppsala, Uppsala, Sweden, 1976)] could provide better definition of between-occasion factor composition, but the number of variables and occasions in this study still surpass the limits of software available.

22. G. Rasch, *Probabilistic Models for Some Intelligence and Attainment Tests* (Danmarks Paedagogiske Institut, Copenhagen, 1960); in *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics* (Univ. of California Press, Berkeley, 1961), vol. 4, pp. 321-33; *Br. J. Math. Stat. Psychol.* 19, 49 (1966).
23. B. Wright and N. Panchapakesan, *Educ. Psychol. Meas.* 29, 23 (1969); B. Wright and R. Mead, *CALFIT: Sample-Free Item Calibration with a Rasch Measurement Model* (Res. Memo. No. 18, Department of Education, Univ. of Chicago, 1975). In using this method, analyses included four blocks of two adjacent measurement points (1 and 2, 2 and 3, 3 and 4, 4 and 5). All test items applied to all children that were common to the two proximate measurement points were included for analysis. Those items that did not fit the theoretical (Rasch) model were not included for further analysis at either measurement point, and what remained were exactly the same items at both points. Between measurement points 1 and 2 there were 126 common items included for analysis; between 2 and 3, 105 items; between 3 and 4, 82 items; and between 4 and 5, 79 items. In no case were there any perfect scores or zero scores.
24. Let M_W = total items after calibration at measurement occasion W ; $C_{W,L+1}$ = items common to both occasion W and occasion $W + 1$; $C_{W,L-1}$ items common to both occasion W and occasion $W - 1$. Since $C_{W,L+1}$ and $C_{W,L-1}$ are subsets of M_W they estimate the same ability. However, a change in origin is necessary to equate the estimates because the computational program centers the scale in an arbitrary origin. Let $X_{W,L-1}$ and $X_{W,L+1}$ be the abilities estimated by using tests of length $C_{W,L-1}$ and $C_{W,L+1}$, respectively. Then:

$$X_{W,L-1} = \beta_0 + \beta X_{W,L+1} \quad (1)$$
 Since the abilities estimated are assumed to be item-free, then the slope in the regression will be equal to 1, and β_0 is the factor by which one ability is shifted to equate with the other. $X_{W,L+1}$ and $X_{W+1,L+1}$ are abilities estimated with one test at two different occasions (note that $C_{W,L+1} = C_{W+1,L-1}$); then by Eq. 1 it is seen that $X_{W,L-1}$ and $X_{W+1,L-1}$ are measuring the same latent trait. Because the scales have different origins, $X_{W+1,L-1}$ is shifted by an amount β_0 to make them comparable.
25. It must be acknowledged here that with this method the interpretability of the data depends upon the comparability of units (ability scores) throughout the range of scores resulting from the Rasch analysis. Although difficult to prove, the argument for equal intervals in the data is strengthened by the fact that the increase in group means prior to treatment is essentially linear. Further discussion of this point may be found in H. Gomez, paper presented at the annual meeting of the American Educational Research Association, New York, 1977.
26. T. W. Anderson, in *Mathematical Methods in the Social Sciences*, K. J. Arrow, S. Karlin, P. Suppes, Eds. (Stanford Univ. Press, Stanford, Calif., 1960).
27. R. D. Bock, *Multivariate Statistical Methods in Behavioral Research* (McGraw-Hill, New York, 1975); in *Problems in Measuring Change*, C. W. Harris, Ed. (Univ. of Wisconsin Press, Madison, 1963), pp. 85-103.
28. Gain during treatment period is defined here as the mean value of a group at a measurement occasion minus the mean value of that group on the previous occasion. Thus the group T1 gains that form the baseline for this analysis are the following: treatment period 1 = .72; period 2 = 1.44; period 3 = 1.74.
29. C. Deutsch, in *Review of Child Development Research*, vol. 3, *Child Development and Social Policy*, B. M. Caldwell and H. N. Ricciuti, Eds. (Univ. of Chicago Press, Chicago, 1973).
30. L. M. Terman and M. A. Merrill, *Stanford-Binet Intelligence Scale, Form L-M* (Houghton, Mifflin, Boston, 1960), adapted for local use.
31. F. Horowitz and L. Paden, in *Review of Child Development Research*, B. M. Caldwell and H. N. Ricciuti, Eds. (Univ. of Chicago Press, Chicago, 1973), vol. 3, pp. 331-335; S. S. Baratz and J. C. Baratz, *Harv. Edu. Rev.* 40, 29 (1970); C. B. Cazden, *Merrill-Palmer Q.* 12, 185 (1966); *Curriculum in Early Childhood Education* (Bernard van Leer Foundation, The Hague, 1974).
32. Colombia has now begun to apply this concept of multiform, integrated attention to its pre-school age children in a nationwide government program in both rural and urban areas. This is, among developing countries, a rarely encountered confluence of science and political decision, and the law creating this social action must be viewed as a very progressive one for Latin

- America. Careful documentation of the results of the program could give additional evidence of the social validity of the scientific findings presented in this article, and could demonstrate the potential value of such programs in the other regions of the world.
33. Adapted from a procedure described by H. G. Birch and A. Lefford, in *Brain Damage in Children: The Biological and Social Aspects*, H. G. Birch, Ed. (Williams & Wilkins, Baltimore, 1964). Only the visual-haptic modality was measured.
 34. The measure was constructed locally using some of the items and format found in C. Bereiter and S. Englemann, *Teaching Disadvantaged Children in the Preschool* (Prentice-Hall, Englewood Cliffs, N. J., 1969), pp. 74–75.
 35. This is a locally modified version of an experimental scale designed by the Growth and Development Unit of the Instituto de Nutrición de Centro América y Panamá, Guatemala.
 36. G. Arthur, *A Point Scale of Performance* (Psychological Corp., New York, 1930). Verbal instructions were developed for the scale and the blocks were enlarged.
 37. D. Wechsler, *WPPSI: Wechsler Preschool and Primary Scale of Intelligence* (Psychological Corp., New York, 1963).
 38. At measurement points 3 and 4, this test is a combination of an adapted version of the arithmetic subscale of the WPPSI and items developed locally. At measurement point 5, the arithmetic test included locally constructed items and an adaptation of the subscale of the WISC-R (39).
 39. D. Wechsler, *WISC-R: Wechsler Intelligence Scale for Children-Revised* (Psychological Corp., New York, 1974).
 40. At measurement points 3 and 4 the mazes test was taken from the WPPSI and at point 5 from the WISC-R.
 41. Taken from (39). The information items in some instances were rewritten because the content was unfamiliar and the order had to be changed when pilot work demonstrated item difficulty levels at variance with the original scale.
 42. At measurement point 4 the test came from the WPPSI, at point 5 from the WISC-R.
 43. At measurement point 4 this was from *WISC: Wechsler Intelligence Scale for Children* (Psychological Corp., New York, 1949); at point 5 the format used was that of the WISC-R.
 44. At measurement point 4 this test was an adaptation from the similarities subscale of the WPPSI. At point 5 it was adapted from the WISC-R. Modifications had to be made similar to those described in (38).
 45. B. Inhelder and J. Piaget, *The Early Growth of Logic in the Child* (Norton, New York, 1964), pp. 151–165. The development of a standardized format for application and scoring was done locally.
 46. This research was supported by grants 700-0634 and 720-0418 of the Ford Foundation and grant 5R01HD07716-02 of the National Institute for Child Health and Human Development. Additional analyses were done under contract No. C-74-0115 of the National Institute of Education. Early financial support was also received from the Medical School of the Universidad del Valle in Cali and the former Council for Intersocietal Studies of Northwestern University, whose members, Lee Sechrest, Donald Campbell, and B. J. Chandler, have provided continual encouragement. Additional financial support from Colombian resources was provided by the Ministerio de Salud, the Ministerio de Educación, and the following private industries: Miles Laboratories de Colombia, Carvajal & Cía., Cementos del Valle, Cartón de Colombia, Colgate-Palmolive de Colombia, La Garantía, and Molinos Pampa Rita.
-