DOCUMENT RESUME

TH AAP INT	TH.	003	762
------------	-----	-----	-----

AUTHOR	Stallings, Jane: Giesen, Phillip A.
TITLE	A Study of Confusability of Codes in Observational
· · ·	Measurement.
INSTITUTION	Stanford Research Inst., Menlo Park, Calif.
PUB DATE	[Apr 74]
NOTE	50p.; Paper presented at the Annual Meeting of the
·	American Educational Research Association (59th
	Chicago, Illinois, April (1974)
-2 ³	
EDRS PRICE	NF-\$0.75 HC-\$3.15 PLUS POSTAGE
DESCRIPTORS	Bias; #Classroom Observation Techniques;
	Codification; Measurement Techniques; *Reliability;
	Simulation; *Video Tape Recordings

ABSTRACT

ED 093 944

Classroom observation has the potential of obtaining valuable information regarding teacher and child behavior. This study examines the accuracy of observers coding a standard stimuli. In this procedure the observer's bias is examined, as well as the confidence that can be placed in the observation code itself. Through these procedures the exact nature of the confusion of codes can be identified. In order to avoid the problems encountered with the paired observer method, an attempt was made to assess the accuracy of observers through the use of controlled videotape examples which allow each interaction (a frame) and sequences of frames to be analyzed for accuracy. Ten videotaped skits were produced to present concise, clear examples of each code used in recording classroom interactions on an observation instrument. Confusability (low observer agreement) matrices were constructed by tallying the observer code sequences. Results of the confusability study identify the specific codes that appear to be reliable as well as those that are confused and need to be redefined. Inter-rater accuracy and videotape simulation accuracy are compared. While the two systems of examining observer accuracy do yield some different information, it is not contradictory, and the videotape system is easier to interpret. (Author/RC)



ED 093944

1

202 762

STANFORD RESEARCH INSTITUTE Menlo Park, California 94025 · U.S.A.

US DEPARTMENT OF HEALTH. EDUCATION & WELFARE NATION& INSTITUTE OF EDUCATION THIS DOCUMENT HAS BEEN REPRO DICED EXACTLY AS RECEIVED FROM THE PERSON OR OFLIANTATION OFFICIN ATING IT POINTS OF VIEW OR OPINIONS STATED DO NOT NEELSSAR, Y REPRE SENT OFFICIAL NATIONAL INSTITUTE OF EDUCATION POSITION OR POLICY

A STUDY OF CONFUSABILITY OF CODES

IN OBSERVATIONAL MEASUREMENT

by

Jane Stallings, Ph.D. Phillip A. Giesen

Presented to the American Educational Research Association 1974 Annual Meeting Chicago, Illinois

April 15-19, 1974

Table of Contents

Intr	oduction . \cdot	1
A.	Procedure to Assess the Confusability of Observation Codes	2
	1. A Description of Procedures	2
	2. Procedural Problems	3
B.	A Description of Confusability Matrices ,	4
c.	Confusion of Observation Codes (A Study Combining the Results of 63 Observers)	11
	1. Findings of "What" Code Confusions	15
	2. Findings for "How" Code Confusions	18
	3. Summery	20
D.	Accuracy of Individual Observers	20
	1. Findings from "What" Codes Occurring Six or More Times	23
•	2. Findings from "How" Codes Occurring Six or More Times	24
	3. Summary	24
E.	A Study Comparing Inter-Rater Accuracy and Videotape Simulation	
, ,		25
	1. Paired Observers	25
10	2. Videotaped Skits (Simulation)	30
	3. Conclusions	36



A STUDY OF CONFUSABILITY OF CODES IN OBSERVATIONAL MEASUREMENT

As a measurement technique, classroom observation has the potential of obtaining valuable information regarding teacher and child behavior. This potential is realized to the extent that the observation data can be shown to be reliable. Reliability, as Cronbach et al (1972) have pointed out, is related to the notion of generalizability from a sample to some universe of interest.

Several sources of unreliability have been identified in past research. Medley & Mitzel (1963) say that:

Most commonly, it [unreliability] occurs when two measures of the same class tend to differ too much; this may happen because the behaviors are unstable, because the observers are unable to agree on what occurs, because the different items which enter into the measurement lack consistency, or for some other reason.

Neither Cronbach nor Medley address the question of the confusability of codes used in observation instruments. The present study examines the accuracy of observers coding a standard stimuli. In this procedure the observer's bias is examined, as well as the confidence that can be placed in the observation code itself. Through these procedures the exact nature of the confusior of codes can be identified.

In previous SRI reliability studies, the technique of pairing the observers with an SRI trainer has been used. However, there are some problems in assessing inter-rater reliability. First, there is some variability in the coding skills of SRI trainers. Second, there is most certainly a variability in the incidents which occur in the classrooms, in what is selected for observation, and in which codes are used in the observations. The optimum arrangement might be to have all observers and SRI trainers observe the same phenomena in the same classroom at the same time. But, as Soar (1973) says:

> The critical problem (of paired observers) is the effect on the classroom of increasing the number of observers. One observer represents a threat to many

teachers and a distraction to the children, at least initially, and as the number of observers increases, these difficulties increase, probably more like a geometric function than an arithmetic one.

In an effort to avoid the problems encountered with the paired observer method, SRI staff has attempted to assess the accuracy of observers through the use of controlled videotape examples. This procedure allows each interaction (or frame) and sequences of frames to be analyzed for accuracy, whereas previously only simple marginal frequency counts of single codes could be computed.

Other investigators in observational research also use videotapes to assess observer accuracy. Soar (1973) used tapes of actual classroom events, and Simmel (1973) cleverly used the last ten minutes of the Johnny Carson Show to Check observer accuracy on a weekly basis. Although they are useful, the limitations of videotapes also should be recognized:

- Because of the difficulty in seeing and hearing, videotapes are more difficult to code than live conversations;
- It is more difficult to understand the gestalt of the situation from a tape than it is from a live situation in the classroom;
- Simulated skits are likely to be more clear-cut examples than those which actually occur in classrooms.

A. PROCEDURE TO ASSESS THE CONFUSABILITY OF OBSERVATION CODES

1. A Description of Procedures

Differing from both Soar and Simmel, SRI staff produced ten videotaped skits. Each simulation is approximately 20 interaction frames long. These skits attempt to present concise, clear examples of each code used in recording classroom interactions on the SRI observation instrument. Each skit begins with a still picture and the voice of a narrator who explains the situation and identifies the focus person. The

These procedures were developed at SRI by J. Philip Baker, Phillip Giesen, and Charles Norwood.



-2-

skit is then shown as regular speed. After the skit is shown once, the still picture and narrator again identify the focus person. Each skit is then shown again, this time with a 2- to 3-second pause between each interaction. The observers are instructed to code this stop-action portion of the skit and to code one frame during each stop or pause.

2. Procedural Problems

3

The reliability coding booklets were returned to SRI and compared with the criterion sequences. This revealed that some observers were coding more than one frame during a pause. Conversely, some observers, possibly while turning pages, omitted frames. The trainers reviewed the coding sequences and deleted extraneous frames or inserted spaces so as to align the observers' sequences with the criterion sequences. Three trainers performed this operation. Since judgment is involved, a check was made on the code sequences of 10 observers to see whether the trainers arranged the sequences in the same manner. The average agreement between trainers in arranging these sequences was 96.4 percent.

Other procedural problems were also encountered due to the experimental nature of the techniques used. Comments received from the observers indicated that not all of the equipment utilized to administer the tapes was in good condition, and, as a result, the sound or pictures were of poor quality. Also, some examples on the criterion tape were technically less than well executed. The most serious problem, however, was that there were too few examples of several of the codes on the criterion tape. Five or fewer examples of a code limited the assurance that representative examples of the code were shown. Further, if an observer missed two out of four possible instances of a code, he only had a score of 50 percent of the criterion correct; however, if he missed two out of 30 possibilities, he had a score of 93 percent of the criterion correct. For this reason, the codes which have fewer than six examples will not be interpreted in this analysis. The number of examples of each code on the complete set of tapes ranges from zero to 40. (This problem is being remedied by the development of more skits.)

-3-

B. A DESCRIPTION OF CONFUSABILITY MATRICES

Confusability of codes refers to codes which were confused with the correct codes by an observer. * Confusability matrices were constructed by tallying the observer code sequences. For each frame, a tally mark was entered in the box or cell created by the juncture of the criterion code and the code marked by the observer. Figure 1 shows an example of a confusability matrix for the "What" codes. The principal diagonal contains the cells indicating correct coding; other cells contain incorrect coding. The column totals are the total number of criterion examples shown on the videotape for each code; the row totals are the total number of times an observer recorded <u>each code</u>, whether correctly or not. An examination of a particular cell reveals whether the code was recorded correctly or incorrectly and, if recorded incorrectly, shows exactly which codes were confused.

The total number of tallies in each cell can be used to calculate the rates of accuracy in two related but distinct ways. The first procedure described above allows an examination of observer bias. If the number in a given cell is compared to the total <u>number of recordings</u> (row total) of the code that pertains to that cell (see the row indicator), a ratio of correct or incorrect responses can be derived. For the cells that fall on the main diagonal, the numbers indicate the proportion of times the code recorded by an observer was correct. For the cells that do not fall on the diagonal, the number indicates the proportion of error. This accuracy rate is called the Accuracy Rate of each observer on each code; i.e., the ratio of correct or incorrect codes of the total number of coded observations.

The point of the second procedure is to assess the confidence that can be placed in each code. The second accuracy figure can be arrived at by comparing the number of tallies located in the same cell to the total number of examples on the criterion tape presented of that specific code (see column total). Again, the number arrived at shows the proportion of correctness to incorrectness, as based on whether the cell falls on the diagonal. This proportion of times the criterion instances were recorded correctly or incorrectly is called the Criterion Accuracy Rate.

* See Table 1 for a brief explanation of the SRI "What" and "How" codes.

-4-

SRI "WHAT" AND "HOW" CODES

"What" Codes	"How" Codes
1 - Command or Request	Н - Нарру
19 - Direct Question	U - Unhappy
2 ~ Open-ended Question	N - Negative
3 - Response	T - Touch
4 - Instruction, Explanation	Q - Question
5 - Comments, Greetings;	G - Guide/Reason
General Action	P - Punish
6 - Task-related Statement	0 - Object
7 - Acknowledge	W - Worth
8 – Praise 9 – Corrective Feedback	DP - Dramatic Play, Pretending
10 - No Response	A - Academic
11 - Waiting	B - Behavior
12 - Observing, Listening	
NV - Nonverbal	
X - Movement	



e .	
	D IN EACH CELL BY ONE OBSERVER
	ONE
	ВΥ
	CELL
	EACH
	II
	RECORDED
	CODES
	"WHAT"
	OF
	TOTAL NUMBER OF "WHAT"
	TOTAL

•

						1) 		l		
Total No. Recorded by Observer	е . С	21	9	28	5	4	00	2	59	17	80	e.	10	
12		1												-
11												·		
10		•												
6													88	
8	6							``				N		
lumn) 7				. T		e.	•			-	80	1		
(by Column) 6 7	 7•	·						J	1	16			•	
								1	9	·	·			
Criterion Examples 4 4NV 5 5NV	·							0						
erion ANV							8					2		
Crit 4						4								
JNV					3	- - -	_							
3				26										-
7			5				2		·			e.		
10		20	-	٦										
1	Э		·		-								-	
"What" Codes	1	10	N .	3	ANE	4	4NV	ي. م	SNV	9	2	8	6	

-6-

ņ

3

n

11

12

26

22

147

. 33

ო

0

00

N

٦

17

~

-

00

~

2

27

9

3 22

Total Number of Examples on Criterion Tape

÷

.

Figure 1

ERIC

Figure 2 presents the proportion of cell tallies to the row totals for each cell. This provides a matrix that presents the Accuracy Rates in place of the raw scores shown on Figure 1. For example, the observer recorded a "1" code correctly three times, so the Accuracy Rate is 100, or 1.00. In the second row, "1Q" was recorded correctly 20 out of 21 times, or 95 percent of the total number of times.

Looking across the "1Q" row, we see that the observer called an example of a "12" a "1Q" five percent of the time. For another example, we look across the row of code "5" and see that our observer did not code any "5's" correctly. Instead, she mistakenly coded two examples of code "5" as a "5NV" and a "6."

In Figure 3, the proportion of criterion examples for each code correctly recorded by the observer are found in the diagonal cells. Entries in cells down the column marked by the correct code other than in the diagonal cells are instances of confusion. As can be seen in the "1Q" column of Figure 3, code "2" and code "3" were sometimes confused with the "1Q" code. The bottom row of the figure presents the correct number of criterion examples of each code which appeared on the videotape. The last column on the table lists the number of times the observer recorded each code. In the code "12" column on Figure 3, the observer recorded three more "12's" than appeared on the videotapes. Apparently, these three were confused with some other code.

Figures 2 and 3 can be overlaid so that the top entry in a cell refers to the accuracy of what the observer recorded, and the lower entry refers to the percent of criterion examples which were correctly coded. Figure 4 illustrates such an overlay. The combined figure tells us that when this observer recorded a "7," it was indeed a "7" (there are no other entries in the code "7" row). However, she only recorded eight "7's" and there were actually 11 examples on the videotapes. Looking down the column for code "7" and at the lower entry in the cell, we see that examples of "7" were recorded as "3" nine percent of the time, as "8" nine percent of the time, and as "12" nine percent of the time. We can conclude from this example that when the observer recorded a "7" it was truly a "7," but she underestimated the number of times' they occurred. She recorded some of the "7's" as "3," "8," or "12."



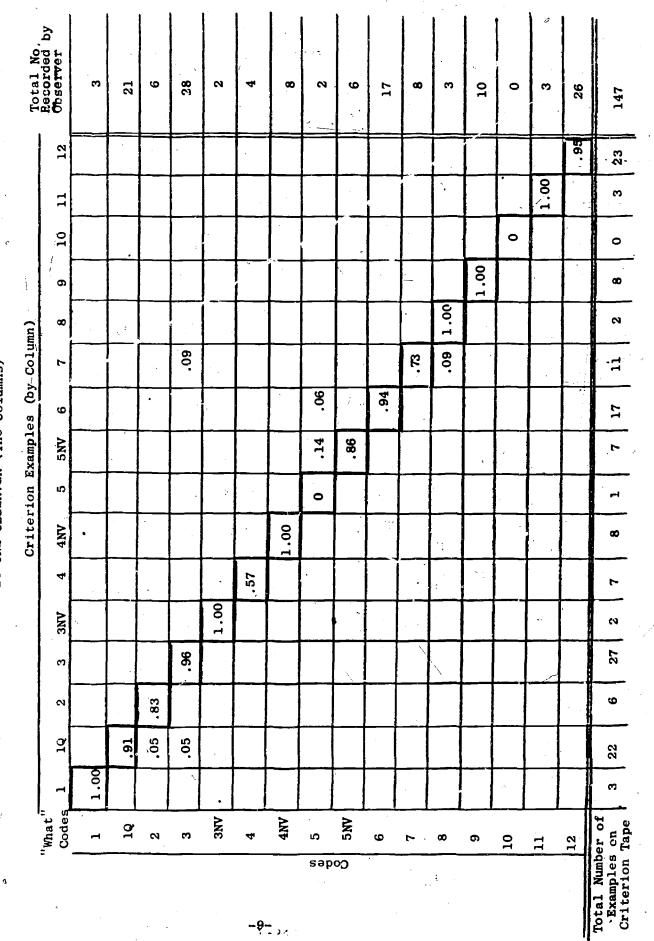
-7-

	What" Codes	n [°]	10	5	თ (ოი	AN BO	() 	19719 4	v sd J	2NA	orde 0	ь 29Я	00 6 82	စ	10	II	12
	ب م	1.00							۰.		23						
	10		.95	.17	.04												
	2			.83							8.	 					
	e	,			.93									.F			
	3NV	1				.1.00	·		·								í.
	4		1	-			1.00										.12
Crit	4NV	, 1						1,00									
Criterion Examples	Ω	·					. ,		, O					. 10			
Examp	5NV								.50	1.00			۱.				
les	9								.50		.94						
	2				.04		х			- ``		1.00	.33		ی ۲		.04
-	. 00							-					.67				
	6					•								.80		· .	
	10	1.0													0		
	п,		-							1						1.00	
	12		.05				_	1									.85
Total No.	Recorded Observer	, ε	21	9	28	2	4	с С	8	9	17	œ	e	10	0	e	26

ERIC Full Text Provided by ERIC

PROPORTIONAL DISTRIBUTION SHOWING HOW CRITERION EXAMPLES WERE RECORDED

BY ONE OBSERVER (The Columns)



ERIC

5		•	Total No.	, Recorded bv Observer		- y ³		g	28	. 0	4	60	2	.9	17	80	n	. 10	0	æ	26	
	·			12		. [.05						at .					` e			. 85 96	
į	ORDED	i		11		•		-					ŀ.			•	2			1.00		
	H REC	ļ.	-	10	к 1		•	а	ģ		;								òò			
	R BOT	•	Celis)	6		•			1			4				L	-	888		;		
•	ER FO		in Cel	00					Ì				ľ				1.00	• •				
	"WHAT" CODES FOR ONE OBSERVER FOR BOTH RECORDED AND CRITERION ACCURACY			2		÷			040			(,	· ·			1. <u>00</u> 73	33				0 0	
•	UT" CODES FOR ONE O CRITERION ACCURACY	•	Examples '(Lower Numbers	9		, s 4						+	50		94 94							
Ф,	S FOR		'(Lowe	SNV			•						50	1.00					.	· · ·		
Figure	, CODE	✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓✓<	ip'les .	2 2	1.	æ .	3						00	1		n e	÷	28				
-1	"WHAT" AND CB	7	Exag	4NV			•.			·		000	•		,							
••	OF 0		Criterion	4							1,00							-	• •		. 12	
λ	PROPORTIONAL DISTRIBUTION	к.	Cri	ANE		*3		•		1.00 1.00								-	· 	· · ·		· · (
	ISTRI	-		ю 1					- 93 96		·			-	, , ,			. 10 .04				
• • •	UAL D		ļ	Ņ				83		·	 				90 17					·····	-	 (
	PORTIC		,	1 0	 		.95 91	.17 .05	04 05							2.						
	PRO			-	1.00	· · ·														{		
- . .		•	, ,	"Viat" Codes	1		 ខ្ព	5	е С	ANE		4NV		SNV		 ·	l	 		<u>1</u>		of
		۳		- 8 - -		(MO			:		xəqwi -10-	N- Je			BGLA	<mark>∂A Op</mark>	<mark>ק pə</mark>	co.cq	<mark>ਿ</mark> ਅਮੁਤਾ	R Beb	<u>Š</u> 12	Total Number of Examples on

•

. -

.

C. CONFUSION OF OBSERVATION CODES (A STUDY COMBINING THE RESULTS OF 63 OBSERVERS)

Analyses of these matrices were used in two ways. First, by combining the results for all observers, the extent of general confusability of codes could be examined. Codes that reveal a high rate of confusion by several observers suggest these possible causes: an overlapping of code definitions, poor videotape eramples, or less-than-adequate training procedures. Second, the accuracy of individual observers could be examined with these matrices. For example, if an observer were not very accurate on code "8" (praise), then codes using "praise" could be examined for anomalies." The findings reported here represent 63 observers spread among 30 geographical locations.

How the observers coded the videotaped examples is shown in Figure 5. The diagonal shows the number of correct codings. The row at the bottom of the table is the number of videotape criterion examples coded by all of the observers. Each figure in the bottom row can be compared with the corresponding cell in the diagonal. For example, code "1" was recorded correctly 160 times out of a possible 245 times. The other entries in the "1" column are sources of confusion.

The proportion of times that the observers were correct in their recordings is presented in Figure 6. For example, the observers recorded "1Q" correctly 80 percent of the time. Looking across the "1Q" row, we see that two percent of the time when a "1Q" was recorded it was truly a "1," and 11 percent of the time it was truly a "2."

The proportion of videotaped examples recorded by the 63 observers is shown in Figure 7. The number in the diagonal reports the percent recorded correctly. The numbers in the columns outside of the diagonal indicate the source of error. If all of the numbers in the columns were in the diagonal cell, the result would be 100 percent correct. The total number of possible examples on the videotape are listed on the bottom row. For example, "1Q" was recorded correctly 77 percent of the time, whereas

Caution: Even if the observer were 100 percent in agreement with the criterion examples, in a study of this type generalization would still be limited by the day-to-day variability of classroom events.

*Sources of error less than three percent are not included on the table.

-11-

, N.	Recorded by Observer	3	3	7	2	0	0	8	υ	6	B							
Tota	Reco Obse	222	1212	31	1447	180	480	548	365	479	853	527	145	439	31	132	1427	8804
	12	12	12	0	19	ß	55	5	9	1	11	9	3	4	5	-	0611	176, 1333
c	11	0	•	0	•	0	, 0	8	0	17	0	0	0	0	М	125	29	176,
	10	0	0	0	0	0		0	0	•	0	0	0	0	Ģ	0	0	0
	6	13	29	e	10	0	œ	0	16	0	29	10	8	378	0	0	4	502
	80	0	0	0	ð	0	0.	2	S	5	0	ъ С	106	1	0	0,	8	121
	7	Ъ	4	ı	14	4	00	13	3	S	36	455	22	2	0	0	37	605
les	9	6	4	0	29	4	7	. 2	200	9	581	F,	` 4	23	0,	0	11	88
Examples	SNV	0	0	0	2	6	1	31	39	343	ຸິ	0	0	0	9	ນີ	19	458
Criterion	- LO	e	0	. 0	4	0	0	0	37	2	÷.	Ō	0	8	0	0	3	56
Crit	4NV	. 0	 0	0	0	0	11	.480	19	9 6	11	0	0	-	0	0	4	622
	.4	e	10	0	3	CV1	344	e l	8	3	21	7	ГТ С	13	1	o .	62	473
	3NV	-	י <mark>הן</mark> הייי	~	17	142		8	-1	e	~	-	0	0	Ħ	, Ģ	9	188
	9 0	ιŋ (24	e C	1310	. 41	34	C.	18	0	75	38	-	6	-		28	1572
	10	5	135.	156	4	0	0	0	0	0	22	-	0	.0.	0	0		320
-	10	13	996	153	35	0	6	-	2	-	31	°.	0	3	e	0	27	1252
	-1	160	27	0	-	8	5	8	14	-	27	ö	0	8		0	2	245* 1252
	"What [™] Codes	, T	10	~~	 	3NN	4	4NV	ى س	SNV		~					8	of over vers
						SIS	VISE	A 05	d beb	000	พ ธย	səp	00					Number of Examples over All Observers
								-12-					- - -		•			N Exai All
			•														•	•

Figure 5

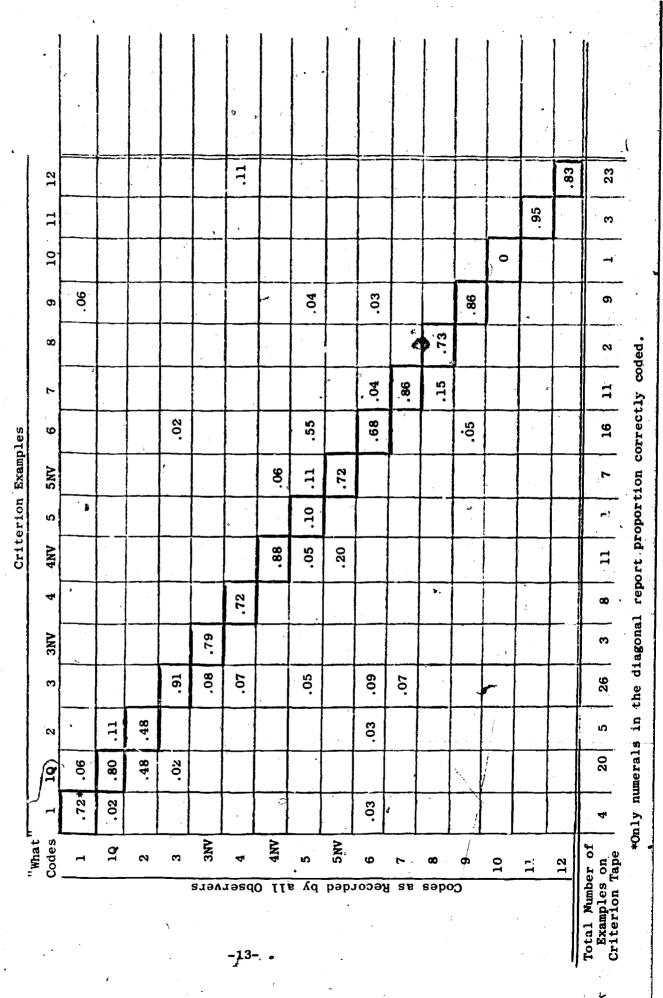
ER

TOTAL OF, "WHAT" CODES RECORDED IN EACH CELL BY 63 OBSERVERS

į a

ERIC

PROPORTIONAL DISTRIBUTION OF CODES RECORDED BY 63 OBSERVERS



"What" , Codes 1 10 2	1 .65* .01	1 q .11 .77 .42	2 .12 .49	3	3NV	4	- 4NV	ი səp	S SNV	6 .11 .07	2	80	5	10	11	12	Total Number of the stamples on the 20 5
3		-	6	.83						7 .05							26
3NV					.76					60.		-					8
4						.73				.04					•	.13	00
4NV							.77		.15							b.	11
5								.66		60.							
5NV									.75								2
5 5NV 6						•		<u> </u>		.66					-		16
2								 		90.	.75				•	8.	• 11
80						<u> </u>			.		.05	.88					N
6	.03	90.				 	 		 	.06			.75				6
, , 10 [°] ,		<u>, , ,</u>							 ,		 			0			
11 1		J.m.				.04		·	.10		 ·	·			.11.	.16 .8	3 23
12	` `		-	-		4										89	6

Figure 7

ŧ

. ;

14

PROPORTIONAL DISTRIBUTION SHOWING CRITERION EXAMPLES RECORDED BY 63 OBSERVERS**

.

4

16 percent of the time it was confused with the codes "1," "2," and "3."

Figure 8 is an overlay of Figures 6 and 7. This provides the data necessary to quickly assess the observers' accuracy (by looking at the topmost entry in a cell) and the percent of criterion codes which have been recorded (by looking at the lower entry in a cell).

1. Findings of "What" Code Confusions

Since a "What" code is required for each interaction, each recorded frame must include a recording of a "What" code. The observer only has the option of recording the correct (or criterion) code or recording the wrong code. The entire frame is considered void if no "What" code is recorded.

In Figure 8, four of the "What" codes have been separated into two categories: the "What" code alone and the "What" code with its "How" modifier. This was done because the meaning or definition of the "What" code is modified or sometimes changed by the addition of these specific "How" codes. An example of this is the "5" code. The definition of the "5" is "general comment," but the definition of the "SNV" is "general action."

As mentioned earlier, the number of criterion examples for some codes is small which limits the conclusions that can be drawn regarding these low frequency codes. For this reason, codes with fewer than six examples on the videotapes will not be discussed.

Nine of the 16 "What" codes have six or more criterion examples of each code. These are the shaded diagonal cells in Figure 8. Of these nine ("1Q," direct question; "3," response; "4," instruction; "4NV," self learning, "5NV," general action; "6," task-related comment; "7," acknowledge; "9," corrective feedback' and "12," attending), only "6" has an observer accuracy rate that is lower than .70.

Code "6," task-related comment, was confused most often with code "3," response. It was also sometimes confused with examples which were actually "1," direct request, "2," open-ended questions, and "9," corrective feedback (see row "6"). This suggests that the definitions and training procedures need to be more exact regarding when to code a taskrelated comment "6." The numbers in the lower section of the cell

-15-

5NV 6 7 8 9 10 11	.03	90.		.02			90.	.11 .55 .04 .04 .09 .23		.04 .04 .06 .06	.05		.05	0			
4 4NV 5	•	· · · · ·		×				.10 .66	.15	04 09						3	
3 3NV			· · · · · ·		8 .79 .76			5		.09 0.09 .0						.13	
0	9 [1	11 42	8.49 2.49		.08	.07		• • • • • • • • • • • • • • • • • • • •		.03 .07							
What Codes 1 10	1 .65 .01	1q .02	2 .12	3 .03	ANE	4	4NV	2	5NV	6 .03 .11	7	60	6	10	11		tal Number of

Figure 8

PROPORTION OF "WHAT" CODES RECORDED IN EACH CELL, CALCULATED OVER 63 OBSERVERS*

ł

*Sources of error less than three percent are not included on this table.

(looking down the "6" column) indicate that 23 percent of the time the criterion examples of "6" were recorded as "5."

The next lowest in reliability was the "4" code, instructing. Eleven percent of the time the observers recorded what was actually "12," observing, as a "4." Since "4" is verbal and "12" is nonverbal, the problem would not appear to be one of confusion in the normal sense but, rather, confusion of which person to focus upon. This conclusion is based on the fact that both of these codes generally occur simultaneously (that is, when a teacher is instructing, "4," the children are usually attending, "12"). Apparently the observers confused which person to record. As can be seen in Figure 8 in the code "12" row, a true example of "4" was sometimes confused and recorded as a "12," which is a further indication that the instructions regarding the focus of observation were not clearly understood by observers.

Code "4NV" describes a child working alone on instructing himself. Observers recorded this reliably 88 percent of the time. They sometimes confused "4NV" with what was truly a "5NV," a code that describes "play" rather than "self instruction in a task." Looking down the "4NV" column, it can be seen that 15 percent of the videotaped examples were recorded as "5NV's." This confusion of "4NV" and "5NV" indicates an overlap of definitions (or a conceptual difficulty in distinguishing "work" from "play").

Criterion examples of code "7," acknowledgement, were sometimes confused with "6" and "12." Code "7" is sometimes confused with code "3," responding (see row 7 in Figure 8). It is easy to see how acknowledging a child can be confused with responding to a child. On the other hand, code "3," responding, was one of the more reliable codes. It was not confused with "7" (see Figure 8). In fact, the observers recorded it correctly 91 percent of the time, and column 3 indicates that five percent of the criterion examples were confused with "6."

Eleven percent of the recorded code "1Q's," asking direct questions, were actually code "2's," asking open-ended questions. The confusion between "1Q" and "2" has long been recognized by the SRI researchers. Each year the variables have been defined more carefully; however, there still seems to be a gray area of unclarity between the two

ERIC[®]

-17-

codes. Code "2," which has too few examples to analyze with confidence, was also confused with "1Q." The results of individual observers were examined, and apparently those observers who observed models which do not often require the "2" code had a higher rate of error.

Eighty-six percent of the time the observers recorded "9," corrective feedback, correctly; five percent of the time code "6" was recorded as "9" (see row 9, the upper value). The criterion examples as illustrated in column 9 (the lower value) were sometimes recorded as "1," "1Q," and "6."

2. Findings for "How" Code Confusions

A "How" code is not always required. This rule leads to four distinct possibilities:

- A required "How" was left out of the frame (omission). These are listed at the bottom of Figure 9;
- A "How" code was recorded when not called for (intrusion). These are listed in the last column of Figure 9;
- The criterion "How" code was confused with another code. These are entered in other than the diagonal cells;
- The criterion "How" code was recorded accurately. These are entered in the diagonal cells.

Caly six of the 14 "How" codes were represented by six or more examples on the videotapes. These are "NV," "X," "A," "B," "DP," and "O" (see Figure 9). As described on page 3, codes with fewer than six examples will not be discussed. Also as described earlier, the upper value in a cell reports the percent of observer accuracy. The lower value in the cell reports the percent of the videotaped examples which were correctly recorded.

The nonverbal code "NV" was recorded correctly 93 percent of the time by observers; and, overall, the observers omitted only 13 percent of the criterion examples. Code "X," movement, was also found to be reasonably reliable. Eighty-nine percent of the time the observer recorded it correctly, but 20 percent of the examples were om ed by observers.

1-1



	T N U H Intrusions	20.	.10	.15	.05	.15	.16	.18		.20 .51	.10	.68	.40 .03 .53 .43	.03 .70 .03 .23 .40	.63 .35	.32 .33 .56 .52	
	ð					 .		}		.13 .66	.87 .75	ł				 9	
(Lower Numbers	Ċ		<u> </u>				<u> </u>									99	1
(Lower	đ	;	 .						0	.13							+
	O							.81 .66								33	┦
on Examples	Ū							├ ──								 39	1
rion	W	·			3		67. 19.									 ļ	╡
Criteri	DP			 		.84 .55									! •	 .43	ļ
	B			-04 13	.95 .43											.43	
	A			.76 .76												.21	
	×		.89											•		.20	
	NV	.93 86							-							 .13	Ī
:	Codes	NN	×	A		pA P		O O T T O T O T O O T O O O T O O O O T O					Z	D	——————————————————————————————————————	Proportion of Omissions	Total Wimbon of

Observers recorded the "A" code, academic, correctly 81 percent of the time. Four percent of the "A's" recorded were truly code "B," and 15 percent of the "A's" were actually intrusions. Seventy-six percent of the videotaped examples were recorded correctly, and 21 percent were omitted.

While 95 percent of the examples recorded as code "B" by the observer were correct (see row B), 43 percent of the "B's" were omitted and 13 percent of the examples of "B's" were incorrectly recorded as "A's." This leads to the conclusion that if a "B" is recorded, it is likely to be correct, but the total number of "B" codes may be underestimated by over 50 percent. An examination of each observer's work is important in order to discover the source of the underestimation. It is possible that only a few observers are grossly underestimating "B's," or it could be that many of the 63 observers are underestimating "B's" to only a small degree.

The two remaining codes with six or more examples ("DP," dramatic play, and "O," use of objects) were recorded accurately over 30 percent of the time, but both codes were underestimated (43 percent s.nd 33 percent of the time).

3. Summary

The results of the confusability study identify the specific codes that appear to be reliable as well as those that are confused and need to be redefined. The findings suggest that some codes, such as "6," "4NV," and "5NV," should be more carefully defined because of overlapping definitions. There is some indication that there should be more careful training of observers on the focus of observation so that "4" and "12" will not be confused. The overall reliability for all observers on the "What" codes was 78 percent and 81 percent for the "How" codes.

D. ACCURACY OF INDIVIDUAL OBSERVERS

The value of this new method for measuring accuracy is that it contributes directly toward interpreting the data. Observer bias can be assessed by examining the overuse, underuse, or confusion of codes. In this study, each observer was responsible for observing one grade level

-20-

at a single site. Therefore, the data collected by each observer is identifiable in the analysis.

In order to determine the accuracy rates for each observer separately, tables were constructed that graphically present, by sponsor, each observer's results (see Table 2). Thus, for example, if an observer in Grade 1 at Site X had difficulty with the code "7," acknowledgment, it is possible to compute the site mean of code "7" and compare it, with the first grade means of code "7" at the four other sites of the sponsor. If the means of the four sites (not in question) are similar and the mean of the site in question differs from the other four, there are two possible explanations: (1) Site X may be truly different from the other four sites, or (2) the observer at Site X may not be recording accurately. In any case, the data resulting from code "7" at Site X would be interpreted with caution. This procedure allows for each observer's data to be reviewed in order to estimate the accuracy of the individual on each code and to allow for the data to be interpreted accordingly.

As an example, Table 2 shows the observer accuracy rate (the top number) and the criterion accuracy rate (the bottom number) for each of the Far West observers for each code. In addition, an overall accuracy rate for each observer on all "What" and "How" codes has been computed and displayed on this table to provide a general idea of the observer's skill." The results are grouped by grade level and site. Similar tables for the other six sponsors in the evaluation were prepared. The complete confusability matrix of all observers is not included in this report but is available at SRI.

As previously discussed on page 3, five or fewer criterion examples of a code minimize the confidence with which the actual results can be

^{*}The overall accuracy rate is arrived at by computing the ratio of correct recordings (those that fall in the diagonal cells) of all codes to the total number of recorded codes and to the total number of criterion instances of the codes. For the "What" codes, the two ratios are the same since the total number of recorded codes is equal to the total number of criterion instances. Two ratios are required for the "How" codes since observers are not required to record a "How" code in each frame which leads to differences between the total numbers of criterion examples and total numbers of recorded codes.



ACCURACY RATES FOR THE "WHAT" AND "HOW" CODES BY FAR WEST OBSERVERS

.59 ,83/ .73 HOH .82/ .85 11/18 5.8 74.7 6.8 61 -38/ OVERALL ACCUTACY RATE THINT. er. .79 .74 18, .78 Ĩ 8. \$ 5 .8. .33/ 0.0/ .33 0.0 8 8 9 <u>.</u>8 1.0/ 1.0/ .75 .50 .50/ 1.0/ 1.0/ .67 .33 .50) 20.1 ò.1 88 રે ક 0.0/1.0/1.0/1.0/1.0/1.0/ 0.0 1.0 .33 1.0 1.0 .50 3 .25/ 0.0/ .33 0.6) 0.1 0.5 .33 .14/ 1.0/ 30, 52 .75/ 0.0/ 1.0/ 1.0 0.0 .35 Þ .38/ . . . 8.5 FIVE OR LESS × .67 0.0/ 1.0/ 0.0 .67 .67 .0. ,0.1 .6 0.0/ 1.0/ 1.0/ 0.0 .50 .50 1.0/ 20.05 -.50/1.0/ 1.0 1.0 .0,1,0/ 1.0,1.0 .35/ 1.0/ 1.0 1.0 0.0/ 1.0/ 0.0 .50 1.0 .50 0.0/1.0/ 0.0 .50 0.0/ 1.0/ è.o.o 0 0.1 0.0 0.0 1.0 .1. .e. 1.0 .67/ 8.8 ð. 9.9) 9 0.1 = CRITERION INSTANCES "HOW" CODES .82/ 1.0/ 1.0 .50 0.0 0.0 .88. 88. 1.0 .88. .88. 88. .13/ 22 .67/1.0/1.0/1.0/ .90 .70 .67 .67 .88/ 1.0/ 1.0/ .54 .78 .25 • 0.0 0.0 .92/ 1.0/ .69 .71 .78/ 88. 16. .88. 1.0 .75/ 8 .31 12.0 . 8 .50/ .50/ .50/ .86/ .46 MORE -3.5 **58**. 18. .66/ .92 **02. .82/ **89**. .89/ .67 8.3 **BIX OR** < .96/1.0/ .93 .83 1.0/ 1.0) o 1 .80 80).1 .0 .0.1 .83. 83 ×).1 .83). 8 1.0/ ે ક .85 **/68**. 96. 8.8 8.2 3.5 ş The accuracy rate is given first and the criterion accuracy rate is given accord: .67/1.0/ .67/ 1.0/ .50/ 1.0/ 1.0 1.0 .67/ 1.0/). . . .67/1.0/ .50/ 1.0/ .50 .67 .43/0.0/0.0/1.0/0.0/ .36/ 0.0/ .25/ 1.0/ 1.0/ .80 0.0 1.0 1.0 1.0 0.0/1.0/1.0/0.0 키 .50 1.0 -1.0 1.0/ 2.8) **0.0** .20/ .50 <u>, 8</u> FIVE OR LESS -.75/ 0.1 /0.1 1.0 /0.1 .80/) 2011 .75/ 8.8 Ň 8.8 . 8 9 9 9 9 §.4 8. 8. 8. 88) 8 9 9 ÷.8 5.8 ~ 200 .1.0 .50 1.0/ 0.1 .13/ 8.8 1.0/ .75/ .05 1.0/ -.93 .75/ .79/ \$2.8 2.8 16. <u>19</u> .83 86. 16.08 92 CRITERION INSTANCES 2 .73/ 1.0/ .80 1.0 .86/ 1.0/ .60 .89 .0.1 /£8. .63 "WHAT" CODES .73/ 2.5 .11/ .88. 88. .75 Φ. કે ક 8.2 .88. .70 70) 0.1 0.2 .90/ .82 90 82 2 ષ્ટ્ર છ .79/ .94 .75/ .69 .7<u>5</u> 85. .47 .33/ 1.0/ .93/ .0,1 È.F. v SIX OR NORE .88/1.0/ .70 .50 .89/ 1.0). 1.0/ 1.0/ 2.2. .86/ .75 14. ∕<u>8</u>8.) 8. 19. **1**.0/ SNV .13 .79/ 1.0 .10 .92/ .83/ .91 1.0/ 82 6.5 1.0/ .90 2N .83/ 18. 1.0/ .11/ .78/ .88 1.0/ .15/ .62/ .89 .15 1.5. .81/ 1.0/ .92/ .96 / 81 ,96. 89. 1.0/ 96 .96. 84 .79/ .92/ .92 80.80 ۳ .94/ .80/ 102. .75/ 95/ 35 84 118. 20: e7 / 2 0209 Salt Lake City 0201 Burkeley 0207 Lebanon 0204 Duluth 0213 Taconi Grade 1 Grado 1 Grade 3 Grado 3 Grade 1 Grade 1 FY SPONSOR Grade 3 Gride 1 Crade 3 Grade 3

ERIC

Accurary Raio: Proportion corruct of the total recorded Criterion Accuracy Raie: Proportion of times the criterion instances were recorded correctly.

Criterion Accuracy Rate

Accuracy Rata

utilized. Therefore, only codes with six or more examples are considered in the analysis of specific grade levels within a site.

As an illustration of how Table 2 can be used, the results of the first observers listed are discussed. The observers are grouped according to the site or project they observed. Going from left to right, the "What" codes are first shown on the extreme left with the codes which are represented by six or more criterion instances. The next section includes the "What" codes that were represented by fewer than six instances. The "How" codes are shown next, with a similar division.

1. Findings from "What" Codes Occurring Six or More Times

The first observer listed, Observer No. 1 from Site A, had an overall reliability rate of .84 on the "What" codes (see Table 2). Of the nine codes with six or more criterion instances, only two codes registered an "observer accuracy" or "criterion accuracy" rate of less than .75. Looking at code "4," instruction, we see an observed accuracy rate of .62 and a criterion accuracy rate of .89. This means that when Observer 1 recorded the "4" code, it was correct 62 percent of the time. The observer actually recorded a "4" 89 percent of the time; thus, she missed only 11 percent of the examples. However, 38 percent of the time when she recorded "4's" she was incorrect. Therefore, variables using the "4" code in the first grade at Site A should be interpreted with caution.

The other code which the observer's results show to be considered less than adequate was the "5NV" code, nonverbal general action. The accuracy rate of 1.00 shows that when she recorded a "5NV" it was always a "5NV"--she did not confuse it. However, she failed to code 50 percent of the videotape examples of "5NV."

The overall results for Observer No. 1 show that the observation data she gathered can be analyzed with a great deal of confidence. Only the "4" and "5NV" code results have to be analyzed with special caution.

Three of the other first grade observers for this sponsor registered accuracy rates of over .70 on the "4" code. The first grade observer at Site E has an accuracy rate of only .54. If the results of the data collection show that Grade 1, Sites A and E, have means and standard deviations for code "4" that differ widely from the other sites,



it may be explained by the observers' confusion in the use of the "4" code.

A similar situation exists with the data for the four other first grade observers on the "5NV" code. The underestimation of the code by Observer No. 1 at Site A is not common to all first grade observers. Therefore, this should be taken into consideration when the data is analyzed.

2. Findings from "How" Codes Occurring Six or More Times

It can be seen on Table 2 that Observer 1 at Site A was 100 percent accurate when she recorded five of the more frequent "How" codes. The one coding exception is "A." academic. Only 67 percent of the time were her "A" recordings correct. Thirty-three percent of the time they were not "A's." However, she recorded 90 percent of the "A's" actually occurring on the videotape. The extra 33 percent that she recorded are considered intrusions,^{*} and they overcestimate the occurrence of this code. Observers at other sites had 'their own specific difficulties, and their data will have to be analyzed in the same way that Observer 1's has been analyzed.

3. Summary

The usefulness of this method of measuring the accuracy of individual observers lies in its capacity to:

- Differentiate codes according to relatively high or low levels of confidence;
- Assess an individual's coding skill on a specific code and examine observer bias;
- Compare individual observer's scores with other observer's scores at the sponsor's same grade level.

By thus identifying the various sources of error in the observation measures, we can more accurately determine whether specific problems lie in the code itself or with the individual observer and interpret the data accordingly.

-24-

See page 18 for an explanation of intrusion.

E. A STUDY COMPARING INTER-RATER ACCURACY AND VIDEOTAPE SIMULATION ACCURACY

The preceding section has examined the confusability of the observation codes and the ability of observers to code criterion videotapes. Videotaped simulation of classroom events are, admittedly, different from actual classroom events. In an effort to compare the accuracy of observer ratings on the simulations and inter-rater accuracy in classrooms, a small study was conducted in one location. This section compares the results obtained from both studies of accuracy for two observers.

1. Paired Observers

The first method, the paired observation, is the most commonly used method of assessing interaction analysis instruments. The procedure followed is to have the two observers situated in the same classroom, coding exactly the same situation simultaneously. The recorded codes are then evaluated in terms of percent agreement between the two observers. Since the speed of the two observers is not expected to be consistent, the ratio of the number of codes recorded by the observer is compared to the ratio of the number of codes recorded by the trainer.

It must be pointed out that this paired observation procedure has some serious limitations. First, two extra people in the classroom are more obtrusive than one. Second, it is almost impossible to assure that the two observers are focusing on exactly the same action. Due to limited space, the two observers may not have the same angle of observation; thus, what they see and hear may be somewhat different and yet each observer could be collecting a correct and adequate sample of the behavior which is occurring. A third problem is that even if the marginal frequency counts of a code by two observers are numerically similar, we cannot be certain that the two observers have recorded specific incidents exactly the same. Similar ratios could occur by chance. Lastly, it happens that during the classroom observations certain interactions or codes do not occur, or occur at such a minimal rate, that reliability cannot be computed.' There is no way to be certain that all codes will be assessed within a given time period.

In the study, data from the 16 five-minute observations were



-25-

examined, using three "Who" codes, twelve "What" codes, and thirteen "How" codes. To assess the coding accuracy of the two observers, the proportion of frames that contained a particular code was recorded for each trainer and trainee. From the proportions, the following equation was computed for each code (p is for the trainer and q is for a given observer on a given code):

The percent agreement = 100 x $\frac{\min (q,p)}{\max (q,p)}^*$

Tables 3 and 4 show the overall percentage reliability of the codes separately in terms of their ratio of frequency. It must be noted that accuracy for low frequency variables is difficult to interpret because if one observer records an event four times and the other only two times and they observe an equal number of frames, the agreement is only 50 percent, even though the actual difference is only two occurrences. Higher frequency variables can tolerate a difference of two occurrences and still show a high percentage of agreement. The data for reach observer is presented separately in Tables 3 and 4. Since there are 16 paired observations, it is possible to have as many as 1,216 frames of interaction. Therefore, we have separated the data into three categories: least frequent, moderately frequent, and most frequent. Table 5 is included to further clarify the results of the paired observations. It includes the frequency scores of the SRI trainer as well as the ratios of occurrence and percent agreement scores for both observers over all codes.

The results show that both of the observers were very reliable on the "Who" codes. The "What" codes were also recorded very reliably, with only two exceptions. Observer 1 recorded less than half as many "8" codes (praise) as the criterion observer, and Observer 2 missed nearly 80 percent of the occurrences of code "6" (task related statement). Significantly, however, both of these codes occurred with low frequency.

The results on the "How" codes were much lower. Observer 1 was quite reliable on the "NV" (nonverbal), "G" (guide to alternative), "A" (academic), and "B" (behavior) codes. She was below the 50 percent

When p = 0 and q = 0, the percent agreement is assigned a value of 100.

-26-

PERCENT AGREEMENT BETWEEN TRAINER/OBSERVER 1

WHO CODES

Percent Agreement 91-100	Least Frequent (0-60)	Moderately Frequent (61-175)	Most Frequent (176-1,216) Adult, Child	Total No. of Codes 2
81-90	Machine	•		1
71-80			•	
61-70				,
51-60	•			
41-50		•	TOTAL	<u> </u>

WHAT CODES

Percent Agreement	Least Frequent (0-60)	Moderately Frequent (61-175)	Most Frequent (176-1,216)	Total No. of Codes
91-100	10		4	2
81-90		1 Q , 6	1, 2, 12	5
71-80	11		· .	1
61-70		5, 9		2
51-60	7		1 1	1
41-50	8		1 ⁹⁰	° 1
			TOTAL	13

HOW CODES

1

Percent Agreement	Least Freque (0-60)	nt Moderately Frequent (61-175)	Most Frequent (176-1,216)	Total No. of Codes	
91-100	U, G, DP		A	4	
81-90					
71-80			NV	1	
61-70		B		1	
51-60					
41-50	х			· 1	
31-40	Т			1	
21-30					
11-20	Q	Н		2	
0-10	N, O, W			3	
	-		TOTAL	13	

-27-

PERCENT AGREEMENT BETWEEN TRAINER/OBSERVER 2

WHO CODES

Percent Agreement 91-100	Least Frequent (0-60)	Moderately Frequent (61-175)	Most Frequent (176-1,216) Adult, Child	Total No. of Codes 2
81–9 0				
71-80				
61-70				
51-60				
41-50	~		τοται	2

WHAT CODES

Percent Agreement	Least Frequent (0-60)	Moderately Frequent (61-175)	Most Frequent (176-1,216)	Total No. of Codes
91-100	8		5	2
81-90			3	1
71-80		4, 7	1Q	· 3
61-70	10		12	2
51–6 0		9	1	2
41-50				
31-40				
21-30			•	
11-20	6	· · · · ·	тута Т.	1
11-20	6		TOTA L	

HOW CODES

Percent Agreement	Least Frequent (0-60)	Moderately Frequent (61-175)	Most Frequent (176-1,216)	Total No. of Codes
91-100	DP		Ā	2
81 -9 0				
71-80			NV (1
61-70				
51-60	Х	•		1
41-50			2	
31-40	Т	*		1
21-30	H, U, Q, G			4
11-20		B		1
0-10	N, O, W		-	3
	• •		TOTAL	13



-28-

PAIRED OBSERVATION RESULTS

	ŕ		0	bserve	r 2			Obser	ver l			
· .			Percent Agreement	Observer Ratio	Trainer Ratio	Traine r Score	• •	Percent Agreement	Observer Ratio	Trainer Ratio	Trainer Score	
		Adult	94	.544	.578	638		97	.482	. 467	502	
		Child	9 8	.431	. 423	467	•	96	.488	.506	544	
	OHM	Machine TOTAL		.026	.000	0*		82	.023	.028	30*	
	3	FREQUEN	NCY	990	1,105	1,105			1,160	1,076	1,076	
		<u> </u>	60	.174	.104	115		88	.171	. 150	161	-
		1Q	73	.106	.145	.160		90	.103	093	100	
		2	*	.001	.000	- 0 *		*	.000	.001	1*	
		3	89	.288	.257	283		88	.212	.242	260	
		4	77	.091	.070	. 77		99 °	.150	.151	162	
		5	92	.091		109		65	.086	056	60	
1		6	17*	.006	.035	38*		81	.069	.056	60	
		7	75*	.039	.052	57*		56*	.020	.036	38*	
	E-	8	100*	.022	.022	24*		44*	.007	.016	17*	
	WHAT	9	55	.038	.069	76 、		65	.042	.065	69	•
	3	10	65*	.031	.020	22*		93*	.013	.014	15*	
		11	*	.020	.000	0*	•	75*	.003	.004	4*	2
		1 2	70	.087	.124	137		83	.121		107	
		NV	71	.144	.203	224		73	.239	.174	187	
		X	51*	.021	.041	45*		43	.053	.023	23*	
		H ¹	22*	.004	.018	20*		11	.006	.057	61	
		U	25*	.016	.004	4*		100*	.000	.000	0*	
	۲,	N	0*	.000	.009	10*		8*	.001	.012	13*	
	ŧ	T	3 3*	.003	.009	9*		33*	.003	.009	9*	
		Q	23*	.005	.022	24*		11*	.001	.009	9*	
		G	29*	.008	.028	31*		100*	.033	.033	35*	
	MOH	Р	*	.002	.000	0*		*	.000	.001	1*	
	Ĥ	0	*	.000	.001	1*		0*	. 000	.016	17*	
		W	*	.000	.002	2*		0*	.000	.007	. 7*	
		DP	100*	.000	.000	0*		100*	.000	.000	· 0*	
		A	93	. 670	.622	687		94	.760	.713	767	
		B	16*	.008	.051	56*		62*	.018	.029	31*	
			-									

* Fewer than 60 criterion instances.

Ratio = occurrence of a specific code/total number of frames recorded



agreement rate for the "X" (movement), "H" (happy), "N" (negative), and "O" (object) codes. The remaining codes occurred less than ten times and, therefore, no accuracy rate could be arrived at.

Observer 2's rate of accuracy was similar on the "How" codes. She was reliable on the "NV," "DP," and "A" codes and below the 50 percent level on the "Q," "G," and "B" codes. Eight of the "How" codes occurred only ten or fewer times, thus generalizations regarding these codes would be made with caution.

2. Videotaped Skits (Simulations)

The second phase of the reliability study was based on videotaped skits. The procedure followed is to have the observers code interactions seen on a videotape and compare that record with predetermined criteria. The tape has stops or pauses between each interaction to insure that each observer knows which interaction to code.

The results are then compiled for each observer, and they reveal both (1) which occurrences were not recorded and (2) which code was erroneously recorded in its place. The procedure allows us to identify the problem codes for each specific observer.

In the figures that follow, two values are shown in each cell. For those cells that fall on the main diagonal, the upper value shows the percent of times the total number of codes recorded was correct. The lower value shows the percent of times the code actually occurred and was recorded correctly by the observer.

For cells that do not fall on the diagonal, the two values indicate proportions of error rather than of accuracy. The upper value shows the percent of times a specific code (as shown by row indicator) was recorded instead of a specific criterion code (indicated by the column) to the total number of recordings of that code. The lower value indicates the percent of times that the specific code (indicated by the row) was recorded when a given criterion code was called for (shown by the column).

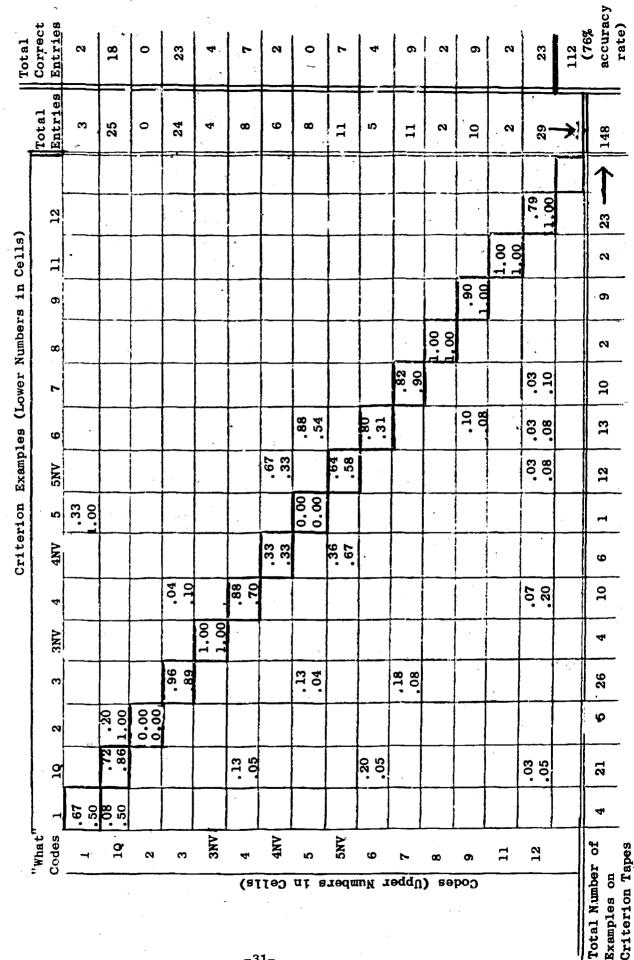
Figures 10, 11, 12, and 13 are matrices showing the percent of accuracy and the percent of the total codes recorded for the two observers. Computations are for the "What" and the "How" codes. The total number of criterion instances of each code is shown at the bottom of each column.



ERIC

PROPORTIONAL DISTRIBUTION OF "WHAT" CODES FOR OBSERVER 1 FOR BOTH RECORDED

AND CRITERION ACCURACY



-31_-

PROPORTIONAL DISTRIBUTION OF "WHAT" CODES FOR OBSERVER 2 FOR BOTH RECORDED

AND CRITERION ACCURACY

(70% accuracy rate) Correct Entries Total . m 100 23 20 N 11 2 2 4 0 0 00 G 4 n ဖ Entries Total 6 13 2 26 2 ŝ ŝ 12 ~ 14 0 e 00 G 30 142 1 . 14 6 19 22 12 . Criterion Examples (Lower Numbers in Cells) 1.00 11 m .22 .11 1.00 .67 φ, ດ 0.00 .25 ŝ 00 10 3 07 £ 09 01. 02. 02. .07 5 .08 20. . 29 .57 .62 5 9 .75 1.00 თ SNV .43 0.00 .30 0.00 .25 .30 ; ۰. 0 ŵ, 8.4 4NV . 10 ٦ .17 **.**80 S 4 .04 1.00 .67 **3NV** ო .89 .82 20 .07 28 .07 ŝ ; 12 80 20 100 100 .07 .20 ŝ 0 85 .44 10 15 .22 .50 .08 25 .14 .25 \$ "What" Codes JNV 4NV SNV 5° 5 12 1 Total number of **Criterion Tapes** က 4 ŝ 9 -Ø 00 (Upper Numbers in (sile) aeboD Examples on

ERIC

-32-

	Correct	_	24	4	34	6	9		4	0	н		L	0		-	87 (78%	accuracy rate)
	Total	Entries	25	5	36	10	7	2	11	1	2	3		1	-	FI I	110	
•	T																	111
suois ston of	sn i q tod o .	T I I	.04	.20	.10	. 10	.14	.50 .05	.64 .33	1.00	.57	.50		1.00			0.00	\uparrow
•	e11s)	H														1.00		8
	in C	Ъ													25		.14	4
	umbers	N				5					.14	j.		0.00	7		.05 .50	2
·	Criterion Examples (Lower Numbers in Cells)	F											1.00				.09 .67	9
URACY	91) 8	a							•		50	50			•			~
ON. ACC	cample	U									• 14 • 00			·				
ITERIC	ion Ey	- d	· ·							88			د			 		0
AND CRITERION ACCURACY	rtter	0							.36 .50	7							. 18	
A		×		·			· <u> </u>	50 60									.05 50	5
		ЪР					86 00	§					 					9
		В				90 75				3.							.14 .25	12
·		A		· · · ·	94 92			• • • • • •									.14 .08	37
		×	•	80 80	-												05	. S
		NV	96 89		i -												14	27
• •	""""""""""""""""""""""""""""""""""""""	Codes	NN	×	¥	<u> </u>	ក្ន (ចារ។	<mark>≥</mark> 70 u	t art	م. مر سمهم	ري 1900]	<u>a</u> U) G	E4 Səpo	י <u>א</u> כי	D	Ħ	Proportion of Omissions	Total Number of Examples on

ERIC FullText Provided by ERIC

·		s Entries	16	ß	27	0	0	0	G	0	0	T	7	0	8		60 (56%	accuracy rate)
5 1 E.	Total	Entrie	17	ß	43	0	0	0	9	. 1 0	5	ч	2	0	2		68	
												·						108
suo fo uot	tenst rodort	u u	.06 .04		.30					1.00 .40	. 50 .04						0.00 0.00	1
RDED	Cells)	Ħ							5							1.00	. 50	8
I RECO	rs in	n	-	·	· · · ·										1.00		.05 .50	4
2 FOR BOTH RECORDED	Number	Z									.50 1.00			8.00 8.00				
2 FOI	ower]	L											1,00 1.00			:		N
"HOW" CODES FOR OBSERVER AND CRITERION ACCUIRACY	Examples (Lower Numbers	0										1.00 .50					. 50	· 0
DR OBS ACCTUR	[dmax2	IJ		-							8:88						+ - +	0
ODES FC	ri on 1	A				-		 .		8:88							┟──┼	
W" CO	Criterion	0							1.00 .86								.02	-
	-	M						8:88									.05 00	
LION C		ad					8:88										.16	
TRIBUT		В			.07	0.00											.23	13
r dis		A	••••••		. 63 77							•					.18	35
LT I ONA		×		1.00 .83													.02	 9
PROPORTIONAL DISTRIBUTION OF		NV	94 59														. 25	27
, 14	"How"			×	A	يم (ع	ີ ຊີ ເເອວ	≥ ur 1	0	0, mN 1	ں nbbe) 89 0	۴ poJ	z	Þ	l	Proportion of Omissions	er of
	ι,					•			-				-			•	Prop of Omi	Total Number Evenules on

The total number that the observer recorded is given at the end of each row.

Those codes that occurred fewer than seven times are listed in the matrices but will not be discussed in the body of this text. A decision was made that, in these cases, the confidence level with which we might make predictions as to the reliability of an observer would be so low as to render it unacceptable. Therefore, only the codes which were tested by seven or more crite ion examples will be considered in this analysis.

As shown in Figure 10, the "What" matrix for Observer 1 indicates that, of the eight codes that included seven or more criterion instances, only the "6" code (task-related statement) and the "5NV" code had a criterion accuracy lower than .70. In the "6" code, 80 percent of the recorded "6" codes were correct, but 69 percent of the criterion codes were missed. Moving up the "6" column we can see that 54 percent of the criterion "6" codes were incorrectly coded as "5." The problem with the "5NV" code is somewhat different. In this case the problem is that both the criterion rate and the observer correctness were low. It appears that on the simulations Observer 1 had difficulty distinguishing the "5NV" code from the "4NV" (self learning or instruction) code, since she often codes the criterion "4NV" instances as "5NV" and also the "5NV" criterion as "4NV."

Over all "What" codes, Observer 1 is reasonably accurate with a criterion rate of .76 which is average for all 63 observers examined by the videotapes on the "What" codes.*

Observer 2 had a reasonable overall criterion accuracy rate (.70) also, but she had coding problems with several codes (see Figure 11). She did not record the "4" code (instruction) 66 percent of the time. The "4NV," "6," "7," and "9" codes were also coded less frequently than required. She used the "12" code (obcerving) eight more times than required. They were confused with codes "3," "4," and "7."

The "How" code accuracy for Observer 1 was also acceptable (see Figure 12). Her overall criterion accuracy rate was .78. This figure indicates that of the 111 criterion "How" codes presented, she recorded them correctly 87 times (see lower right hand corner of Figure 12).

This figure is computed by dividing the total number of correct entries by the exact number of videotaped criterion examples.

-35-

The only "How" code that Observer 1 recorded with less than 70 percent accuracy was the "O" code (objects); 50 percent were missed and 64 percent were recorded when not indicated. The other codes that fell below a .70 rate of accuracy were codes that included fewer than seven criterion instances.

Observer 2 had a more difficult time recording the "How" codes from the tapes. In Figure 13, her overall accuracy rate is shown as only .56. On individual codes, the "O" was very reliable (1.09/.86), but the "NV" was not used 41 percent of the time required. The "A" (academic) was coded when not called for sixteen times as well as omitted eight times when it should have been coded. "B" (behavior) and "DP" (dramatic play) were ignored completely.

3. Conclusions

Two distinct procedures, the videotaped skits and the paired observations, were used to assess the accuracy of two observers. The results indicate average reliability for both observers on the "What" code category. For the "How" category, Observer 1 is above average, but Observer 2 is below the average of the other 62 observers.

Specifically, Observer 1 was acceptably accurate on the more frequently used individual codes. Many of her codes, such as the "1Q," "3," "4," "7," "9," "12," "NV," "A," and "B" were shown to be very reliable on both procedures. Only the "O" code (use of objects) was shown to be unreliable on both procedures.

The results were equally good for Observer 2 on the "What" codes with only the "6" code (task related comment) being shown unreliable on both procedures. The "How" code "B" (behavior) was also recorded poorly in both procedures. In the case of the videotape codings she missed the 13 examples of the "B" code and underestimated it in the paired observations. On the "A" code (academic), Observer 2 was 93 percent accurate on the paired observations but had a .63/.77 reliability on the videotapes. Other "How" codes such as "movement" and "object" are acceptably accurate on the videotapes while "nonverbal" is acceptably accurate on the paired observations. "Guide" and "question," which were underestimated in the inter-rater analysis, have too few examples on the videotape to be discussed in terms of reliability.

While simulated videotaped events are limited in their scope



-36-

and differ from the classroom situation, they do offer a standard stimulus to examine each observer's ability to code specified events and to identify observer bias. There is still some confounding in the source of "system error"; however, the variation introduced by a second observer is eliminated. While the two systems of examining observer accuracy do yield some different information, it is not contradictory, and the videotape system is by fat asier to interpret.