

16 / A44 NCI Program
Scanning and Compression Project
Working Paper No. 4
September 18, 1973

EVALUATION OF
IMAGE COMPRESSION FOR
INSURANCE INDUSTRY DOCUMENTS



IBM IBM IBM IBM IB

F. B. Wood
SDD Los Gatos
622-5265

R. B. Arps
SDD Los Gatos
622-5687

Confidential

Until December 31, 1978
For IBM Internal Use thereafter

Evaluation of Image Compression for
Insurance Industry Documents

ABSTRACT

A "representative set" and a "additional set" of insurance industry documents were scanned and digitized in Fochester with a nominal 120 X 120 pel/inch resolution with 64 levels of gray scale. Their digital tape records were thresholded to black/white in Los Gatos.

The documents were then run through a series of black pel counting, prediction, run-length counting, and analysis programs to obtain statistics on:

- (1) Compressed page distributions and document percent black vs. compressed record size (for market planners)
- (2) Compressed scan line distributions and the ratio between low and high resolution storage requirements (for system architects)
- (3) The impact on compression performance of changing dual base counters (for image product developers).

F. B. Wood
F. B. Arps

SUMMARY

The basic NCI algorithm used in the NCI Prototype compresses representative 180 kilo-byte documents to a range of 9 to 60 kilo-bytes with an average of around 30 kilo-bytes. A "pathologic" document that doesn't compress was found to require 181 Kilo-bytes of compressed storage. It has been examined and was found to be a copy on heat sensitive copier paper that had darkened with age.

The percent blackness for all the documents were plotted against their respective compressed record size. The points clustered into groups of documents of the same class, such as "applications", "claim reports", "change memos", etc. Further the clusters fell close to a slope line of 5 kilo-bytes per percent black. This rule of thumb of 5 kilo-bytes for each percent of blackness appears to hold for documents up to 15% blackness.

The compressed scan line analysis including both representative and additional documents shows that 1056-pel scan lines average to 176-bits per compressed line. 27.7% compress to 4-bits, 70.9% compress to values in the range 8 to 1056 bits, and 1.4% expand to the range 1060 to 2004 bits.

A subset of these documents was scanned at both 120 X 120 pel/inch and 240 x 240 pel/inch resolutions and analyzed. For equivalent threshold conditions, a document scanned at 240 X 240 pel/inch requires 2.6 times the compressed storage needed when it is scanned at 120 X 120 pel/inch.

Testing different dual-base counters with the 3-pel NCI predictor shows that programming different run length counters for separate classes of documents can give 7% more compression over using a fixed (10,6) counter.

TABLE OF CONTENTS

- I. Description of the Documents
 - A. Classification
 - B. Examples of the Representative Set
 - C. "Pathologic Case" Example
 - D. Typical Compressed Size for Different Insurance Forms

Table I. Representative Set of Documents Compressed by NCI Algorithms
- II. Sample Distribution of Compressed Page sizes for NCI Algorithm
 - Table II. Representative Set Compression Results for NCI.
 - Table III. Additional Set Compression Results for NCI Algorithm
- III. Relation Between Original Percent Black and Resultant Compressed Image Size for NCI Algorithm
- IV. Distribution of Compressed Scan Line Sizes for NCI Algorithm
- V. Comparison of Low & High Resolution Documents
 - Table IV. Comparison of Low and High Resolution Los Gatos Raster Scanned Representative Set
- VI. NCI Optimizing Algorithm
 - Table V. Representative Set plus Additional Set Summary of Optimizing Algorithm Results
- VII. References and Notes
 - Appendix A. Description of the NCI Algorithm
 - Appendix B. List of Scanned Documents

List of Figures

- Fig. 1. "Application for State Farm Insurance",
Doc. #1.0A
- Fig. 2. "Change of Address", Doc. #7.1A
- Fig. 3. "Drivers License File" (printout), Doc. #3.0A
- Fig. 4. "Auto Underwriters Memo", (faded copy), Doc. #7.7B
- Fig. 5. Sample Distribution of Compressed Document size
for Representative Set using NCI Algorithm
(2 kilo-byte intervals)
- Fig. 6. Sample distribution of Compressed Document size
for Representative plus Additional Set
using NCI Algorithm. (2 Kilo-byte intervals)
- Fig. 7. "Percent Black vs. Compressed Image Size"
- Fig. 8. Sample Distribution for 92752 Compressed Scan Lines
(68 documents) Full Scale to Show Maximum
- Fig. 9. Data for Fig. 8 at larger scale.
L=4 point is truncated.
Curve from geometric model is also included.
- Fig. 10. Expanded scale to show the few lines that expanded.
- Fig. 11. Relationship between Compressed Record Size for
High and Low Resolution Scanned Images of the
Representative Set.
- Fig. 12. Compressed Document Size vs. Dual Base Code for
Representative Set with NCI Predictor for all 2-bit,
3-bit, 4-bit, and 5-bit Dual Base Counters
- Fig. 13. Predictor Table and Run Length Counter Coding
for NCI Algorithm

I. DESCRIPTION OF DOCUMENTS

A. Classification

Sixty-eight documents were scanned at nominal⁽¹⁾ 120 by 120 pels/inch resolution. The documents were numbered to indicate the class of document as shown below:

1. Applications
2. Inspection Reports
3. Motor Vehicle Reports (MVP's)
4. Renewal Notices which contain some form of correspondence from the policyholder
5. Policyholder correspondence
6. Change Requests
7. Correspondence from Agents
8. Automobile Claim Reports
9. Copies of Claims Drafts paid by Field Claims representatives
10. Polaroid pictures of damaged cars

The complete list of documents and copies of a few samples are included in Working Paper No. 3⁽²⁾. Three of the 21 representative documents are reproduced here, and the "pathologic" one from the set of 47 additional documents is reproduced here.

B. Examples of "Representative Set" Documents

Figure 1:

The "Application for State Farm Automobile Insurance" represents the greatest detail and lowest compression ratio for a full page in the set of representative documents.

"Application For..." #1.0A (48)
Compression Ratio = 3.02 to 1.
Compressed from 1,440,384 to 476,972 bits
(59,622 bytes)

Figure 2:

The "Change of Address" form, when compressed individually comes close to having the average compression ratio as determined collectively for the

NCI Algorithm in the representative set of documents
(8-1/2" x 11" field)

"Change of Address" #7.1A(61)
Compression Ratio = 6.90 to 1.
Compressed from 1,440,384 to 208,648 bits
(26,081 bytes)

Figure 3:

The computer printout (5" x 3-1/2" data in an 8-1/2" x 11" field) of drivers license file data represents the document with the largest compression ratio within the representative set.

"A303101 13M....." #3.0A (52)
Compression Ratio = 20.42 to 1.
Compressed from 1,440,384 bits to 70,541 bits
(8,818 bytes)

C. "Pathologic Case" Example:

Figure 4:

The "Auto Underwriting Memo" which was scanned from a heat sensitive thermofax copy had sufficient dark background that it can be considered a pathologic case. This document had 21.5% black, while other "application" forms had between 11% and 15% black. The background of the faded thermofax paper added sufficient noise so that the document failed to compress. It took 1% more bits than the original to store the "compressed" output.

"Auto Underwriting Memo" #7.7B (45)
Compression Ratio = 0.99 to 1.
The compression algorithm expanded this document
from 1,440,384 bits to 1,452,300 bits.

D. Typical Compressed Size for Different Insurance Forms

The "representative set" of documents, together with their insurance form numbers, are listed in Table I. The compressed sizes obtained with the NCI algorithm are included in the last column. The compressed size ranges between 8.9 kilo-bytes and 60.0 kilo-bytes for this set.

AG 4089 IL.2 APPLICATION FOR STATE FARM AUTOMOBILE INSURANCE OFFICE COPY

*Plus
my*

6788 691 13

VEH. 1	NEW	REINS.	TRANS	ADDED	QUALIFYING POLICY NO.	CLASS	REPLACES POLICY NO.	Office Use
VEH. 2			X				2743 895-E21-46 D.	
NAME PLEASE PRINT LAST NAME: RUIZ FIRST NAME: TARZAN MIDDLE NAME OR INITIAL: F.								DEC 11 1972
MAILING ADDRESS NUMBER AND STREET: PO Box 128 CITY: Hull STATE: Ill ZIP CODE: 62343								
RESIDENCE If other than Mailing Address NUMBER AND STREET: 8 mi East of Hannibal, Mo on Rt. 36. CITY: Hull STATE: Ill ZIP CODE: 62343								
EXACT LOCATION If residence, if no street number used FORMER ADDRESS If in community less than 6 months 8119 Duiven St. - Norfolk Va.								
DURING THE PAST 5 YEARS, HAS THE APPLICANT OR ANY HOUSEHOLD MEMBER HAD AUTO INSURANCE CANCELLED, BEEN REFUSED ISSUANCE OR RENEWAL, OR RECEIVED NOTICE OF SUCH INTENT? Yes <input type="checkbox"/> No <input checked="" type="checkbox"/>								
DURING THE PAST 5 YEARS, HAS THE APPLICANT, ANY HOUSEHOLD MEMBER, OR ANY REGULAR DRIVER: a. Had license to drive or registration suspended, revoked or refused? Yes <input type="checkbox"/> No <input checked="" type="checkbox"/> b. Been the driver in any automobile accident or loss? Yes <input checked="" type="checkbox"/> No <input type="checkbox"/> c. Been convicted for forfeited bail for traffic violations? Yes <input type="checkbox"/> No <input checked="" type="checkbox"/>								
DOES THE APPLICANT OR ANY REGULAR DRIVER HAVE ANY: 1. PHYSICAL LIMITATIONS OR 2. MENTAL DEFECTS? IF YES, EXPLAIN IN REMARKS. Yes <input type="checkbox"/> No <input checked="" type="checkbox"/>								
IF YES, WAS IT WITHIN THE PAST YEAR? Yes <input type="checkbox"/> No <input checked="" type="checkbox"/>								
MOST RECENT LIABILITY INSURER: SF COMPANY-EXPLAIN IF NONE: POLICY NO.: FROM: MONTH-DAY-YEAR TO: MONTH-DAY-YEAR								
VEH. 1	YEAR	MAKE AND MODEL	BODY TYPE	CYLS.	VEHICLE IDENTIFICATION NO.			
	72	FORD - GALAXIE	4 dr sed.	8	2N54H126471			
VEH. 2	YEAR	MAKE AND MODEL	BODY TYPE	CYLS.	VEHICLE IDENTIFICATION NO.			
AIR CONDITIONING YES <input type="checkbox"/> NO <input checked="" type="checkbox"/> MO. AND YR. PURCHASED 9/72 NEW USED <input checked="" type="checkbox"/> COST INCL. EQUIP. \$4900 AMT. OWED \$8000 OTHER INSURANCE ENGINE DISPLACEMENT FOR HIGH PERFORMANCE VEH. ONLY HORSE-POWER UTILITY VEH. COST NEW								
LIENHOLDER: Bank of Virginia MAILING ADDRESS: 21st & Granby St. - Norfolk Va. ZIP CODE: 23517								
NO. CARS IN HOUSEHOLD: 1 NO. INS. BY STATE FARM: 1 EXISTING DAMAGE OR MODIFIED: 1 VEH. 1 NO YES-DESCRIBE: 1 VEH. 2 NO YES-DESCRIBE: 2								
RATING: VEH. 1 TERR. NO. IN CITY BUSINESS USE FARM RATE DRIVEN TO AND FROM WORK MILES ONE WAY AVERAGE WEEKLY ESTIMATED ANNUAL MILEAGE 2 CAR RULE DRIVER TRAINING GOOD STUDENT CHILDREN AT SCHOOL RAYING GROUP CLASS ACC. GRP.								
1 40 No - - 0 8000-10,000 - - - - - 7 180 -								
COVERAGES: The insurance applied for is only for the coverages indicated by specific premium entry. If premium cannot be entered, check boxes to indicate coverage requested.								
BIPD 25/50/15 5000/10/25 MED. PAY. LIMITS 5000 COMPRE-HENSIVE 0 80% COLLISION DEDUCTIBLE AMOUNTS 50 COLLISION EMERGENCY ROAD SERVICE RENTAL REIMBURSEMENT UNINSURED MOTOR VEHICLE 10/20 Limits Same as BI 10/20 Limits Same as BI AUTO DEATH INDEMNITY SPECIFIC DISABILITY LOSS OF EARNINGS ENDORSEMENTS								
S & Z SECTION: FULL NAME OF PERSONS TO BE INSURED: S-PRINCIPAL SUM AMOUNT PREM. Z PREMIUM								
MPP ACCOUNT NO. AGENT'S CODE STAMP: DATE AND TIME OF APPLICATION MO. DAY YR. 12 7 72 AM PM 3 3 HUNSAKER 27								
TOTALS: 90.90								
REMARKS:								
SIGNATURES: BINDER EFFECTIVE DATE: 12-7-72 STATE FARM MUTUAL AUTOMOBILE INSURANCE COMPANY STATE FARM FIRE AND CASUALTY COMPANY AGENT'S SIGNATURE: <i>Milton E. Bredler</i> APPLICANT'S SIGNATURE: <i>Tarzan F. Ruiz</i>								

Fig. 1. "Application for State Farm Insurance", Doc. #1.0A

CHANGE OF ADDRESS

IMPORTANT: If any policy has a different named insured, indicate the name immediately after the policy number.

NAME OF INSURED: George Thomas Current

NEW ADDRESS: PERMANENT TEMPORARY, HOW LONG? _____

RESIDENCE: 315 E. Baird St. Olney Ill
NUMBER STREET CITY STATE ZIP CODE COUNTY TOWNSHIP

MAILING: _____
NUMBER STREET CITY STATE ZIP CODE

EFFECTIVE DATE: 12-6-72 MPP ACCOUNT NO: _____

16
 DEC 08 1972

<input checked="" type="checkbox"/> AUTO	POLICY NUMBERS	CLASS CHANGES		POSTDATE/REAGENT	FARM RATED		TERR. NO.
		FROM	TO		YES	NO	
	<u>6595-903-F09BD</u>				<input type="checkbox"/>	<input type="checkbox"/>	
	<u>6651 543-114-13C</u>				<input type="checkbox"/>	<input type="checkbox"/>	
					<input type="checkbox"/>	<input type="checkbox"/>	
					<input type="checkbox"/>	<input type="checkbox"/>	INSIDE CITY LIMITS? <input type="checkbox"/> YES <input type="checkbox"/> NO

LIFE POLICY NUMBERS

FIRE INSURED MOVED: INTO WITHIN OUT OF MY TERRITORY

POLICY NUMBERS	CANCEL EFFECTIVE	REWRITE POLICY (ATTACH APP)	ENDORSE REMOVAL (COMPLETE)
_____	_____	<input type="checkbox"/>	<input type="checkbox"/>
_____	_____	<input type="checkbox"/>	<input type="checkbox"/>
_____	_____	<input type="checkbox"/>	<input type="checkbox"/>

CONSTRUCTION	YEAR BUILT	TYPE OF ROOF	OCCUPIED AS	NUMBER OF FAMILIES	BOARDS OR ROOMERS
PROTECTION CLASS	DISTANCE TO HYDRANT	FEET	TO FIRE DEPT. MILES	IF OUTSIDE CITY, DISTANCE FROM CITY LIMITS	MILES
NAME OF SERVICING FIRE DEPARTMENT		NAME OF FIRE DISTRICT (WHERE APPLICABLE)		ANNUAL CONTENTS RATES AT NEW LOCATION	
NUMBER SQ. FEET	FIRST FLOOR	SECOND FLOOR	OCCUPIED BY	OWNER	TENANT
			VACANT	IF VACANT, EXPLAIN	

(1) FOR NON-RESIDENTIAL CONTENTS REMOVAL, USE APPLICATION F7-865 (2) ATTACH SNAPSHOT (AS REQUIRED BY REGIONAL OFFICE)

HEALTH POLICY NUMBERS

DATE: 12-6-72 AGENT'S NAME AND CODE: R O McDaniel 8456

MS 5941

Fig. 2. "Change of Address", Doc. #7.1A

A303161 13 M 72-12-05 12-07-72 L200-6521-2200 00620

THE FOLLOWING INFORMATION IS FURNISHED FROM THE DRIVERS
LICENSE FILE OF THE PERSON IDENTIFIED ABOVE PURSUANT TO THE
PROVISIONS OF THE ILLINOIS VEHICLE CODE.

CTIS L LLCY
800 MARKLAND L200-6521-2200
SALEM 62881 BIRTH DATE

SEX HT. WT. HAIR EYES TYPE ISS DATE CLASS RESCRIC EXP DATE
M 68 205 BRWN BRWN 1 08 09 71 A* 100 07 14 74

TYPE	DATE OF	DATE OF	DESC OF	ACC CR	TERM	DATE	STOP IN
ACTION	ARREST	ACTION	ACTION	DCCKET NO.	CF	SLSP	EFFECT
99	092468	093068	107000	0001706	00000		00
YRCW							
99	041770	041770	1061000	0116994	00000		00

FTC

Fig. 3. "Drivers License File" (printout), Doc. #3.0A

(E)

AUTO UNDERWRITING MEMO



10-7-72
 FROM: Mr. J. J. [unclear] 1-1-72
 POLICY NUMBERS: 159991130 NAME: Francis Robert [unclear]
1599104613 ADDRESS: [unclear]

PLEASE REPLY TO ITEMS CIRCLED

- 1. Contact _____ and secure details of accident(s)/violation(s) in which he/she was involved on _____.
- 2. The change of name recently requested will require a fully completed and signed application.
- 3. Obtain a release of credits from _____.
- 4. Forward a driver training/good student certificate for _____.
- 5. _____ has a health problem. In order to properly review the risk we need to see that a medical report be reviewed by our Medical Department. Please have him and/or her visit the COMPLETE NAMES AND ADDRESS of the doctor who have treated him and/or her, and/or where he and/or she has been confined recently, and the enclosed medical authorization forms. If he and/or she has not been treated recently, see the last doctor he and/or she visited. After completion, have the forms signed and return them to us. (If individual is a minor have parents sign.)
- 6. Furnish driving directions to insured's residence, plus references.
- 7. If a driver has moved, Obtain current address and zip code.
- 8. _____ of the insured during the past year and length of time at each.

- 9. The vehicle insured under this policy was a loss on _____. What disposition is to be made of this policy?
- 10. Forward drivers license number, date issued and date of birth for _____.
- 11. The driver license number for _____ is incorrect. The State has no record for the license given. Forward correct one.
- 12. The insured reportedly is not a citizen of the U.S. Does he/she have a permanent visa form? Will he/she be in the U.S. 12 months?
- 13. Driver information on _____
 Full Name _____
 Date of Birth _____
 Driver License No. _____
 Occupation _____
 Prior Conviction(s) _____
 Any physical defects _____
 Dates and details of past accidents and violations _____
- 14. Review indicated class _____ because _____
115 [unclear] [unclear]
115 [unclear] [unclear]

REPLY BELOW 2nd and REQUEST - PLEASE REPLY DC 20 1972

insured was to drop off [unclear] information
have [unclear] will call insured again
the week to get the information
I will follow up in one week to make
sure information obtained

DATE 12/4/72 SIGNATURE Stencoro 17

Fig. 4. "Auto Underwriters Memo", (faded copy), Doc. #7.7B

TABLE I

State Farm Insurance
 Representative Set of Documents
 Compressed by NCI Algorithm
 (120 X 120 pel/inch Resolution)

<u>Document Number</u>	<u>Insurance Form No.</u>	<u>Description</u>	<u>Compressed Size (Kilo Bytes)</u>
1.0A*	AG 4069 IL.2	Application "Ruiz"	60.0
1.0B*	AG 4069 IL.2	Bottom of Application "Ruiz"	18.3
2.0A*	6033	Automobile Insurance Report "Weil"	55.0
2.0B*	6033-F	Backside of Auto Insurance Report "Weil"	28.8
3.0A		Computer Printout from MV License File	8.9
4.0A*	G 54864.1	Renewal Notice "Smith"	12.1
4.0B*	G 54864.1	Backside of Renewal Notice	11.4
5.0A*	G 54844.2	Follow-up Renewal Notice "Kroll"	13.8
5.0B		Handwritten letter "Kroll"	17.4
5.1A		Handwritten letter "Mathes"	34.5
5.1B		Second page of handwritten letter	19.5
6.0A	G 4178.2	Change Request "Vlahan"	44.6
7.0A*	AG 4709.12	Change Memo "Belz"	33.3
7.1A	MG 5941	Change of Address "Current"	26.2
7.2A		Policy holder Review "Burns"	20.3
7.3A	N 4371.4	Inter-Office Correspondence "Keeran"	13.5

<u>Document Number</u>	<u>Insurance Form No.</u>	<u>Description</u>	<u>Compressed Size (Kilo Bytes)</u>
7.4A	AG 4070	Binder for SF MVI "Dunes"	28.4
7.5A*	MG 4517.6a	Agents MEMO (change) "Pomack"	39.4
8.0A*	G 4684n	Automobile Claim Report "Clark"	55.7
9.0A*		Claim Draft "Mesbit"	18.0

* This document and others of the same type have been plotted on the sheet "Per Cent Black vs. Compressed Image Size" to show clustering. (Fig. 7)

II. SAMPLE DISTRIBUTION FOR NCI ALGORITHM

Each of the 21 documents from the "representative set" and each of the 47 documents in the "additional set" were run through the compression calculating programs for the NCI Algorithm. The results are tabulated in Table II and III. Those tables are for images obtained by scanning on the PIDAS scanner in IPM Rochester and computer thresholding at Los Gatos. The document numbers are the same as in Working Paper No. 3.

The second column in Table II is the compression ratio obtained with the NCI algorithm. The "bound" shown in parentheses is the theoretical limits of compression for that document with the NCI predictor. The fourth column "Efficiency" gives the efficiency of the (10,6) dual base run length codes when used with the predictor output for the given document. The fifth and sixth columns ("Compressor Kilo Bits/Kilo-bytes) give the compressed image size in kilo-bits with the size in kilo-bytes in parentheses.

The distribution of compressed documents sizes are plotted in Fig. 5.

Gaussian curves for the same mean and variance are also plotted for comparison.

The compressed image sizes are plotted in Figure 6 for the NCI algorithm for both the "representative set" and the "addition Set". A bar has been added to show the image size of the raw data. The "pathologic case" is also included.

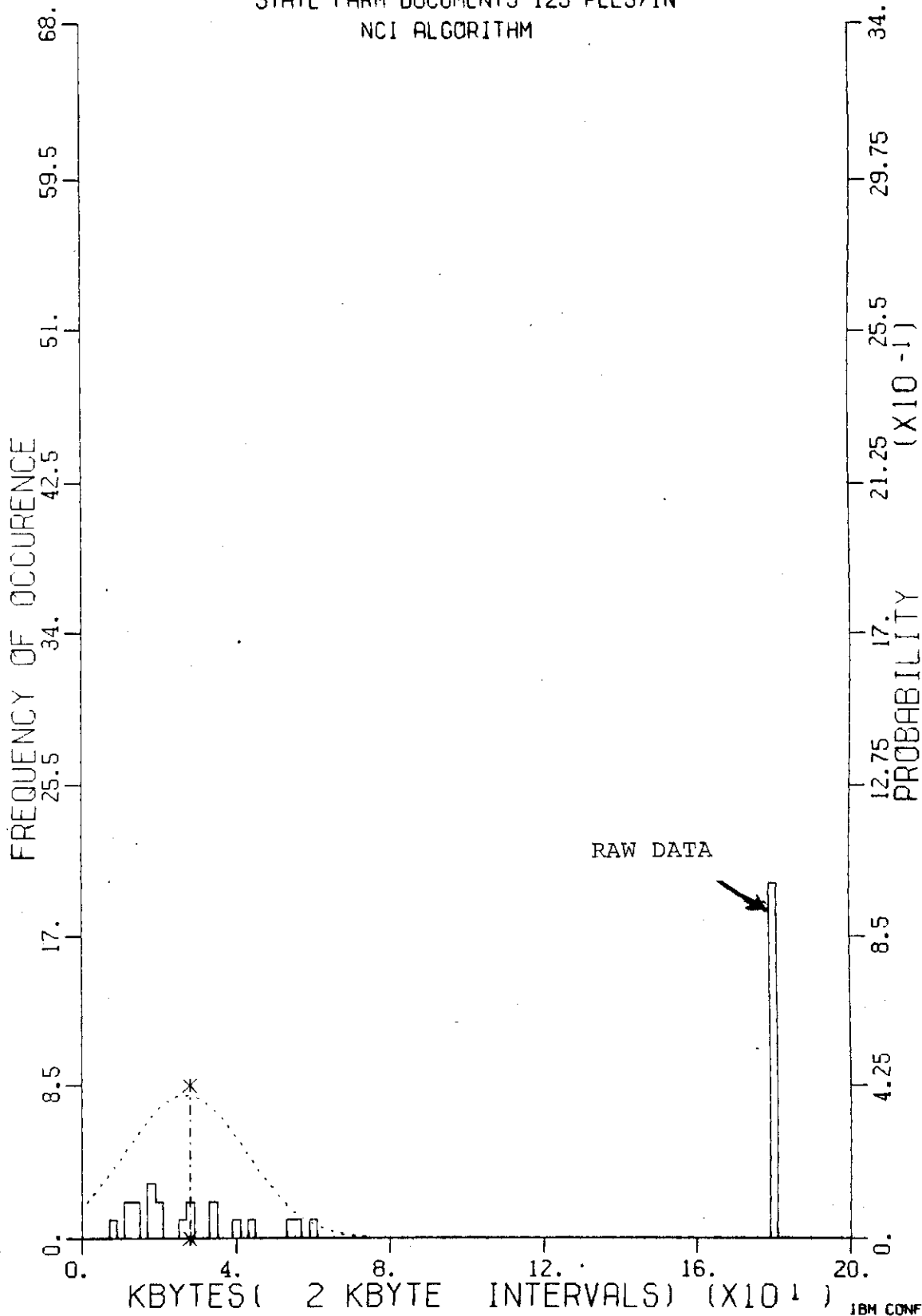
TABLE II

STATE FARM DOCUMENTS
 Representative Set of 21 Documents
 1,440,384 kilo bits
 (180 kilo-bytes)

Doc. Num.	SDDSIO	NCT ALGORITHM		Compressed Kbits (Kbytes)
		Comp. (bound)	Eff. %	
1.0A	48	3.0 (3.6)	85.0	477 (59.6)
1.0B	49	9.9 (11.2)	87.8	146 (18.3)
2.0A	50	3.3 (4.0)	82.7	440 (55.0)
2.0B	51	6.3 (7.4)	84.5	230 (28.8)
3.0A	52	20.4 (23.9)	85.6	71 (8.9)
4.0A	53	14.8 (17.5)	84.5	98 (12.2)
4.0B	54	15.8 (18.3)	87.8	91 (11.4)
5.0A	55	13.0 (15.3)	85.3	110 (13.8)
5.0B	56	10.3 (11.5)	90.1	139 (17.4)
5.1A	57	5.2 (5.8)	89.6	276 (34.5)
5.1B	58	9.3 (10.4)	89.2	156 (19.5)
6.0A	59	4.0 (5.7)	70.7	357 (44.6)
7.0A	60	5.4 (6.3)	85.3	266 (33.2)
7.1A	61	6.9 (7.9)	87.0	209 (26.2)
7.2A	62	8.8 (10.0)	88.8	162 (20.3)
7.3A	63	13.4 (15.0)	88.8	108 (13.5)
7.4A	64	6.4 (7.2)	88.4	227 (28.4)
7.5A	65	4.6 (5.4)	84.9	314 (39.2)
8.0A	66	3.2 (3.8)	85.5	445 (55.6)
9.0A	67	10.0 (11.6)	86.3	144 (18.0)
10.0A		13.6 (16.9)	80.4	107 (13.4)

Min. 3.0A 71 Kbits 8.88 Kbytes
 Max. 1.0A 477 59.7

SAMPLE DISTRIBUTION FOR 20 COMPRESSED DOCUMENTS
STATE FARM DOCUMENTS 125 PELS/IN
NCI ALGORITHM



IBM CONFIDENTIAL

Fig. 5. Sample Distribution of Compressed Document size for Representative Set using NCI Algorithm (2 kilo-byte intervals)

SAMPLE DISTRIBUTION FOR 88 COMPRESSED DOCUMENTS
STATE FARM DOCUMENTS 125 PELS/IN
NCI ALGORITHM

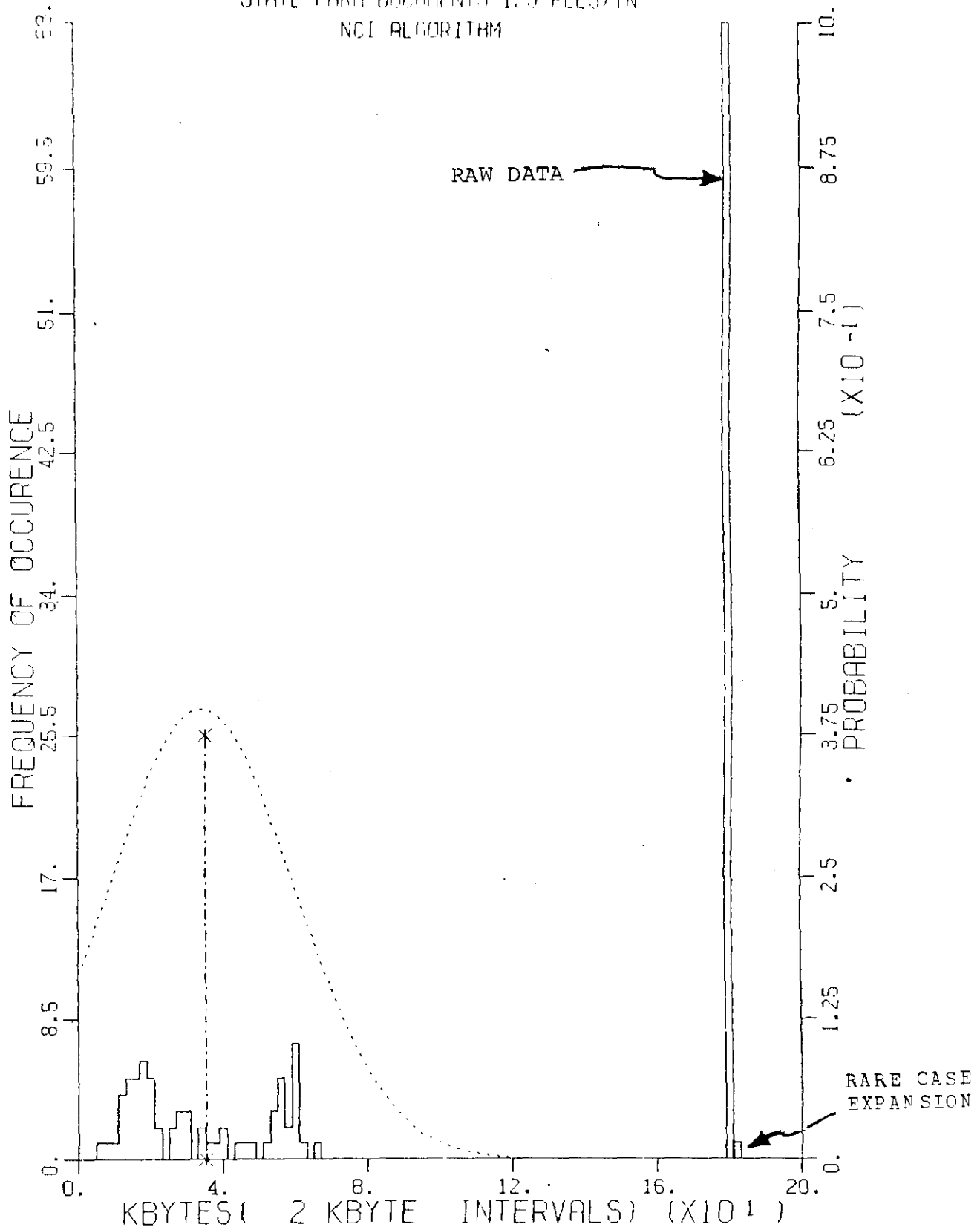


Fig. 6. Sample distribution of Compressed Document size for Representative plus Additional Set using NCI Algorithm. (2 Kilo-byte intervals)

TABLE III

STATE FARM DOCUMENTS
 Additional Set of 47 Documents

Doc. Num.	NCI ALGORITHM		
	Comp. (Bound)	Eff. %	Compressed K-Bits
1.1A	3.2 (3.8)	85.5	446
1.1B	11.6 (13.2)	87.9	124
1.2A	3.2 (3.8)	85.4	448
1.2B	13.3 (15.1)	87.7	109
1.3A	2.0 (3.5)	85.1	487
1.3B	10.6 (12.1)	87.7	136
1.4A	3.0 (3.5)	85.3	484
1.4B	11.0 (12.5)	88.0	131
1.5A	3.2 (3.7)	85.8	454
1.5B	10.5 (11.8)	88.5	137
1.6A	3.0 (3.5)	85.8	485
1.6B	14.6 (16.8)	87.0	96
1.7A	3.0 (3.5)	85.5	479
1.7B	9.2 (10.5)	87.5	157
1.8A	3.3 (3.9)	85.0	440
1.8B	27.1 (32.0)	84.6	53
1.9A	2.8 (3.2)	85.2	524
1.9B	9.7 (11.1)	87.6	149
2.1A	3.4 (4.1)	82.1	429
2.1B	6.8 (8.2)	83.7	211
2.2A	3.1 (3.8)	83.3	460
2.2B	5.1 (5.9)	86.0	282
2.3A	2.9 (3.5)	83.1	502
2.3B	4.7 (5.4)	86.0	309
6.1A	6.6 (7.8)	85.1	218
6.1B	9.3 (11.0)	84.9	154
6.1C	11.3 (13.1)	86.2	128
8.1A	3.3 (3.8)	86.0	438
8.2A	3.8 (4.6)	83.9	377
8.3A	3.1 (3.6)	84.9	473
8.4A	3.4 (4.0)	85.9	422
8.5A	3.0 (3.5)	86.1	483
8.6A	3.1 (3.6)	85.2	470
5.1A	18.0 (21.3)	84.5	80

NCI ALGORITHM

<u>Doc.</u> <u>Num.</u>	<u>Comp. (Bound)</u>	<u>Eff. %</u>	<u>Compressed</u> <u>K-Bits</u>
5.1B	10.0 (11.4)	87.6	145
5.2A	6.0 (6.5)	91.6	241
5.2B	8.2 (9.4)	87.8	175
5.2C	9.0 (10.2)	88.1	160
5.2D	8.3 (9.8)	84.3	174
9.1A	11.3 (13.2)	85.3	128
9.2A	12.7 (15.1)	84.5	113
7.6A	6.1 (7.0)	86.9	238
7.6B	16.0 (18.9)	84.9	90
7.7A	6.0 (7.1)	83.5	242
7.7B	** 0.99 (1.40)	70.7	1,452
7.8A	4.6 (5.4)	86.0	313
7.8B	4.0 (4.6)	86.1	365

** Image that expanded

III. RELATION BETWEEN PERCENT BLACK IN DOCUMENT AND THE RESULTANT COMPRESSED IMAGE SIZE

The complete set of 68 documents were run through two computer programs, one program counts the number of black pels in the black/white source document. The other uses the NCI Algorithm predictor and dual base (10,6) run length codes to calculate the compressed image size.

The percent black was then plotted versus the compressed image size in Figure 7. Only the documents for which there were two or more of the same type are included in the figure. It was apparent that documents of similar type clustered about particular regions on the percent black vs. compressed image size plot. In figure 7 the clusters are labelled with the class of document and the insurance industry form numbers.

The dashed line represents a slope of 5 kilo-bytes of compressed image per 1% black in a 180 kilo-bytes source document.

IV. DISTRIBUTION OF COMPRESSED SCAN LINE SIZE FOR NCI ALGORITHM

All documents in the "representative set" and the "additional set", a total of 68 documents, were compressed using Frank Nemeč's assembly language implementation of the NCI algorithm. The distribution of compressed scan line lengths for 68 x 1364 = 92752 lines is plotted in Figure 8. The frequency distribution for different compressed line lengths ran between 1 and 25732, which is too large a range to show in a linear scale. Therefore, three different scales are used in Figs. 8, 9, and 10. 27.7% of the scan lines compress to 4-bits and 70.9% compress to values in the range of 8 to 1056 bits. The remaining 1.4% expand to the range 1060 to 2004 bits.

The scale resolution has been increased in Fig. 9 to show more detail of the distribution and to compare the actual distribution with an empirically fitted geometric mass distribution function, $P(J)$:

$$P(J) = 0.277 \quad \text{for } J=1$$
$$= 0.011(0.986)^{(J-2)} \quad \text{for } J=2,3,\dots,501$$

where J is the number of compressed half-bytes,
and $L(\text{bits}) = 4(J)$ (half-bytes)

The scale of Fig. 9 is still too small to see the detail of the few lines that expanded to almost twice input size. The same data is plotted in Fig. 10 for a finer scale to show the distribution of the 1.4% of the lines that expanded.

V. COMPARISON OF LOW AND HIGH RESOLUTION DOCUMENTS

The "representative set" of documents were rapidly scanned by the Los Gatos Raster Scanner at 120 x 120 pel/inch and 240 x 240 pel/inch resolution under equivalent threshold conditions to quickly obtain a comparison of the same documents for low and high resolution.

A complete analysis of Table IV indicates that the source black pels for high resolution is slightly

PER CENT BLACK VS. COMPRESSED IMAGE SIZE
 STATE FARM DOCUMENTS 125 PELS/IN

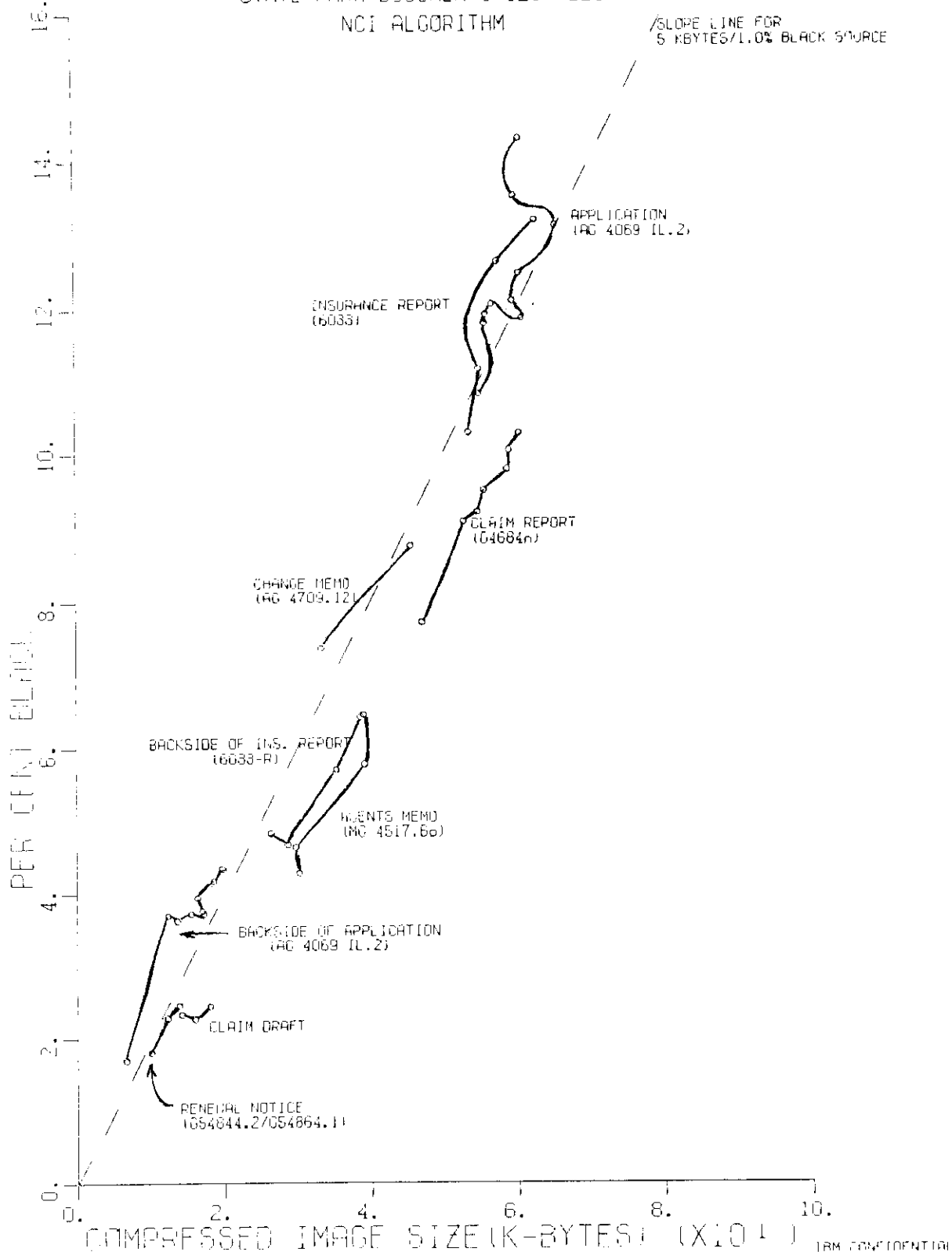


Fig. 7. "Percent Black vs. Compressed Image Size"

P.E. 4445
F.B. 4000
F.H. NEMEL
P.M. PLTAR
08/29/78
68 0005

SAMPLE DISTRIBUTION FOR 92752 COMPRESSED SCAN LINES
STATE FARM DOCUMENTS 125 PELS/INCH
NCT ALGORITHM

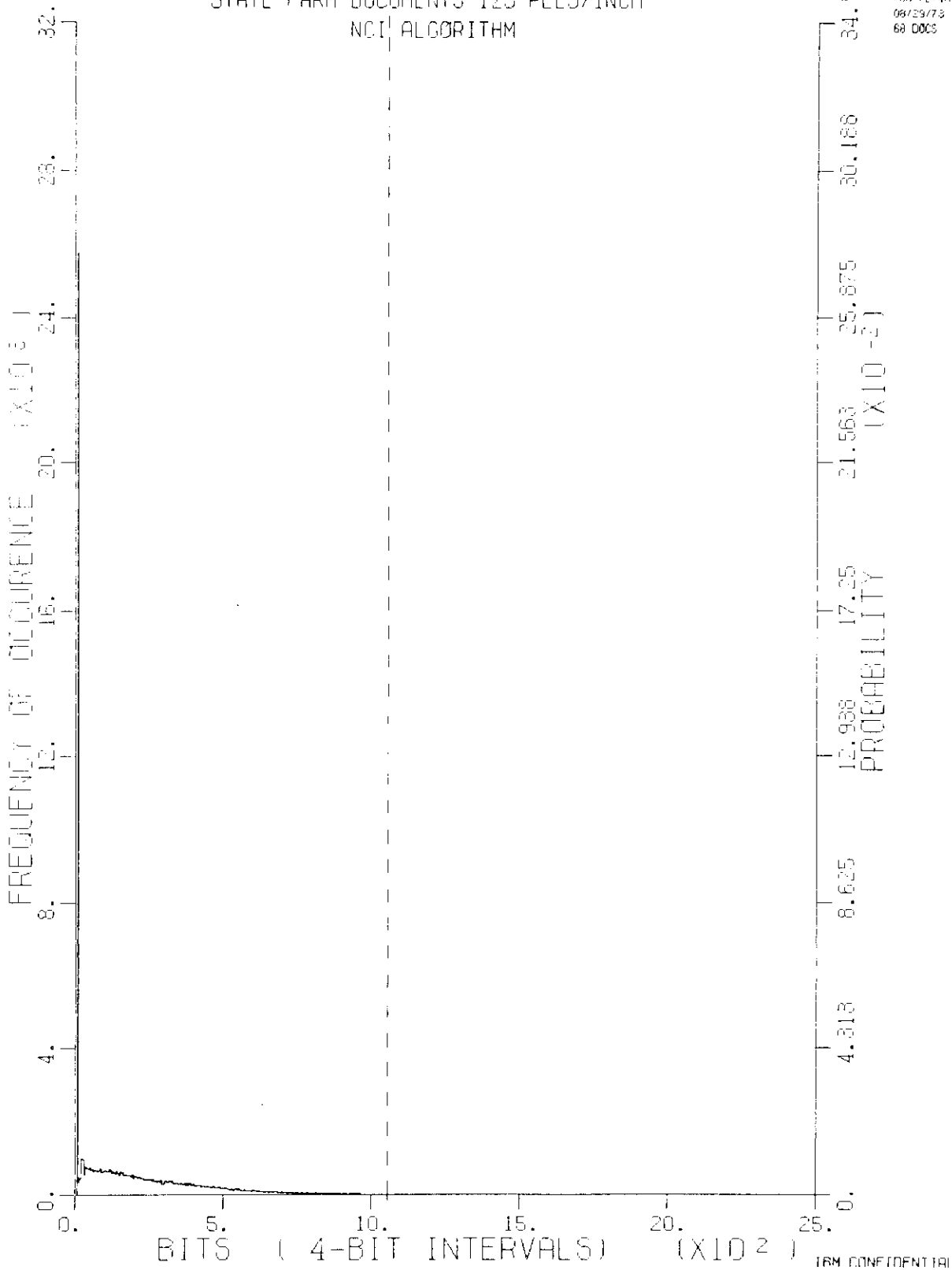


Fig. 8. Sample Distribution for 92752 Compressed Scan Lines (68 documents) Full Scale to Show Maximum

IBM CONFIDENTIAL

SAMPLE DISTRIBUTION FOR 92752 COMPRESSED SCAN LINES
 STATE FARM DOCUMENTS 120 PELS/INCH
 NCI ALGORITHM

P.L. HOFF
 F.B. WOOD
 F.W. NEMEL
 P.M. PELI
 09/13/73
 88 0005
 1125 PELS

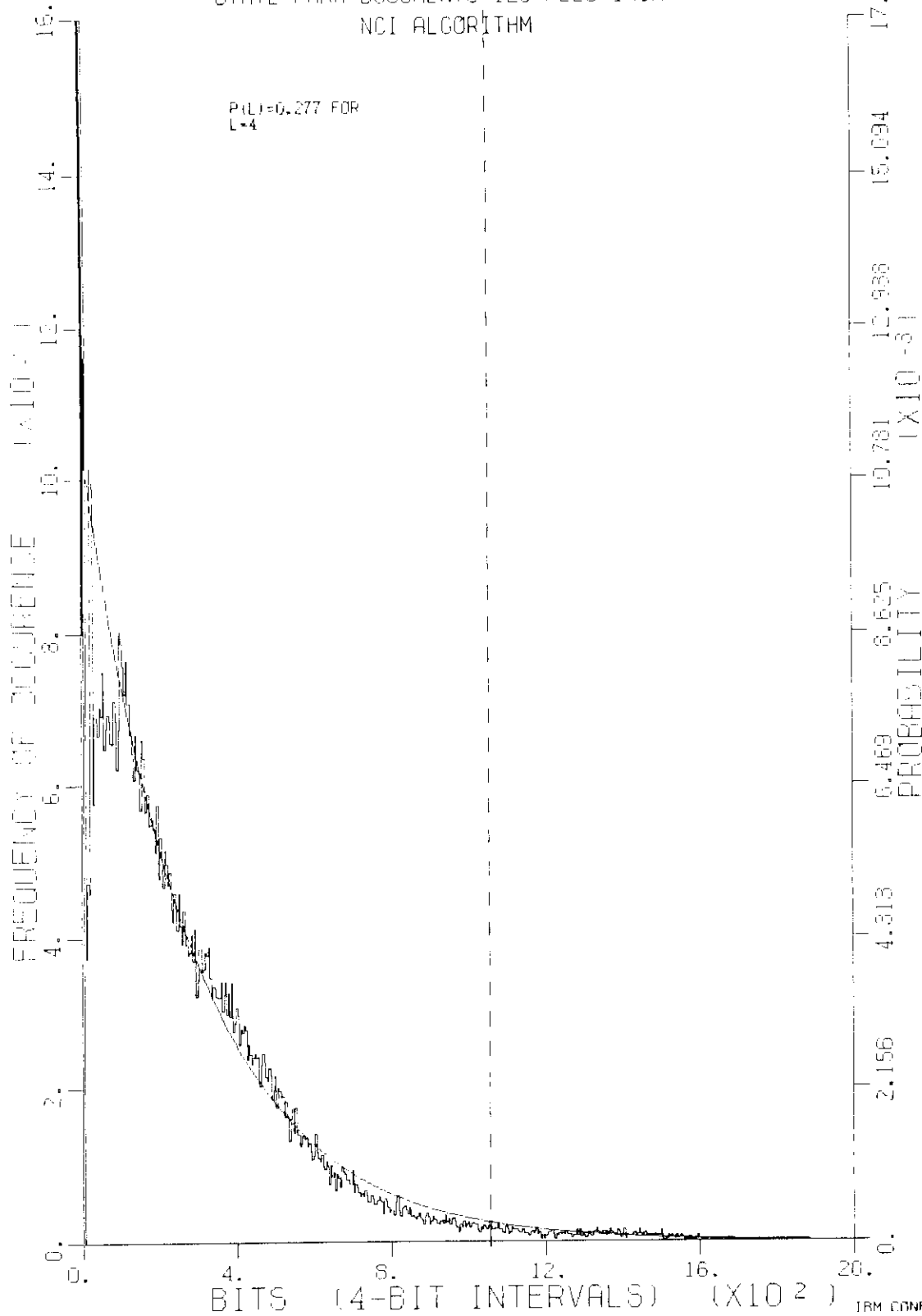


Fig. 9. Data for Fig. 8 at larger scale.
 L=4 point is truncated.
 Curve from geometric model is also included.

SAMPLE DISTRIBUTION FOR 92752 COMPRESSED SCAN LINES
STATE FARM DOCUMENTS 125 PELS/INCH
NOI ALGORITHM

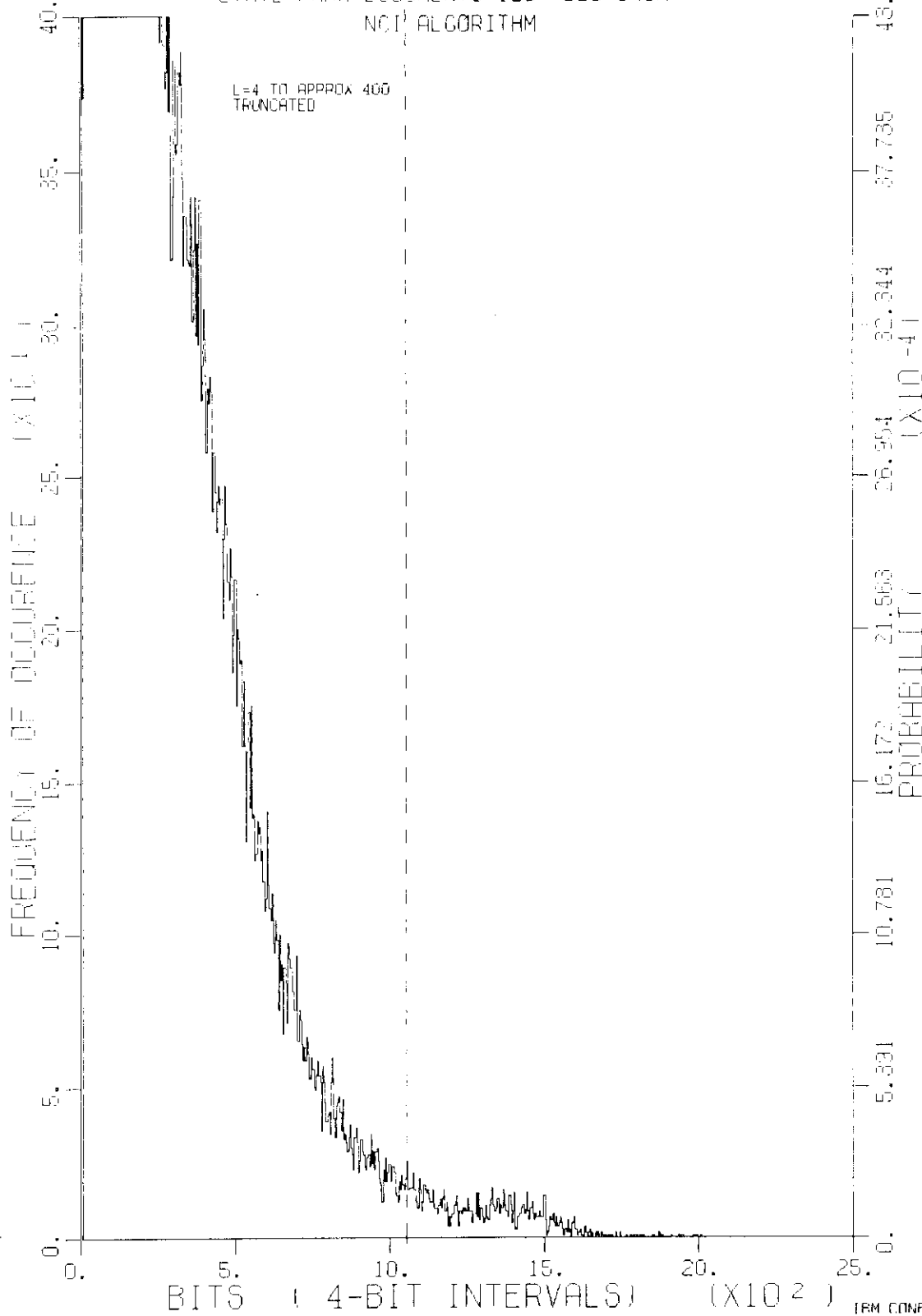


Fig. 10. Expanded scale to show the few lines that expanded.

higher than four times the black pels at low resolution. This is reasonably close to the theoretical factor of four that would be anticipated from doubling the resolution in both dimensions.

The compressed record sizes at high resolution are plotted against the compressed record sizes at low resolution in Fig. 11. A line indicating a 2.6X ratio between low and high resolution compressed size has been added as a rule of thumb. (Note that the theoretical lower bound is a ratio of 2.3X).

This data is slightly different than the PIDAS scanned, dynamic thresholded data used in the rest of this report. However, the ratio should be very indicative of the algorithm's resolution response.

VI. NCI OPTIMIZED ALGORITHM

The 21 documents in the "representative set" were concatenated as one large document of dimension 8 1/2" X 231" in order to explore optimal dual base counting schemes and obtain maximum coder efficiency. This concatenated document was run through the NCI 3-bit predictor followed by a program which computed the compressed document size for all possible dual base run length coders having 2-bit, 3-bit, 4-bit, or 5-bit counters.

The dual base counters follow the formula:

$$\text{Count} = a_0 (p^0 n^0) + a_1 (p^1 n^0) + a_2 (p^1 n^1) \\ + a_3 (p^1 n^2) + \dots + a_i (p^1 n^{i-1})$$

For 2-bit counters $p + n = 2^2 = 4$, giving $2^2 - 1 = 3$ possible dual base counter systems:

$$(p, n) = (1, 3), (2, 2), (3, 1)$$

For 3-bit counter $p + n = 2^3 = 8$; $2^3 - 1 = 7$ and the possible dual base counting schemes are:

$$(p, n) = (1, 7), (2, 6), (3, 5), (4, 4), (5, 3), (6, 2), (7, 1)$$

The results are plotted in Figure 12. The compression bound for this concatenated data set for the NCI predictor is 7.693. The largest compression ratio was obtained for the (5,3) - code, giving a compressed image size of 4,251,882 bits with a compression ratio of 7.114, which is an efficiency of 92.5%.

A number of other codes had fairly good efficiency. They are listed in Table V.

data is padded out to an integral byte boundary, if necessary.

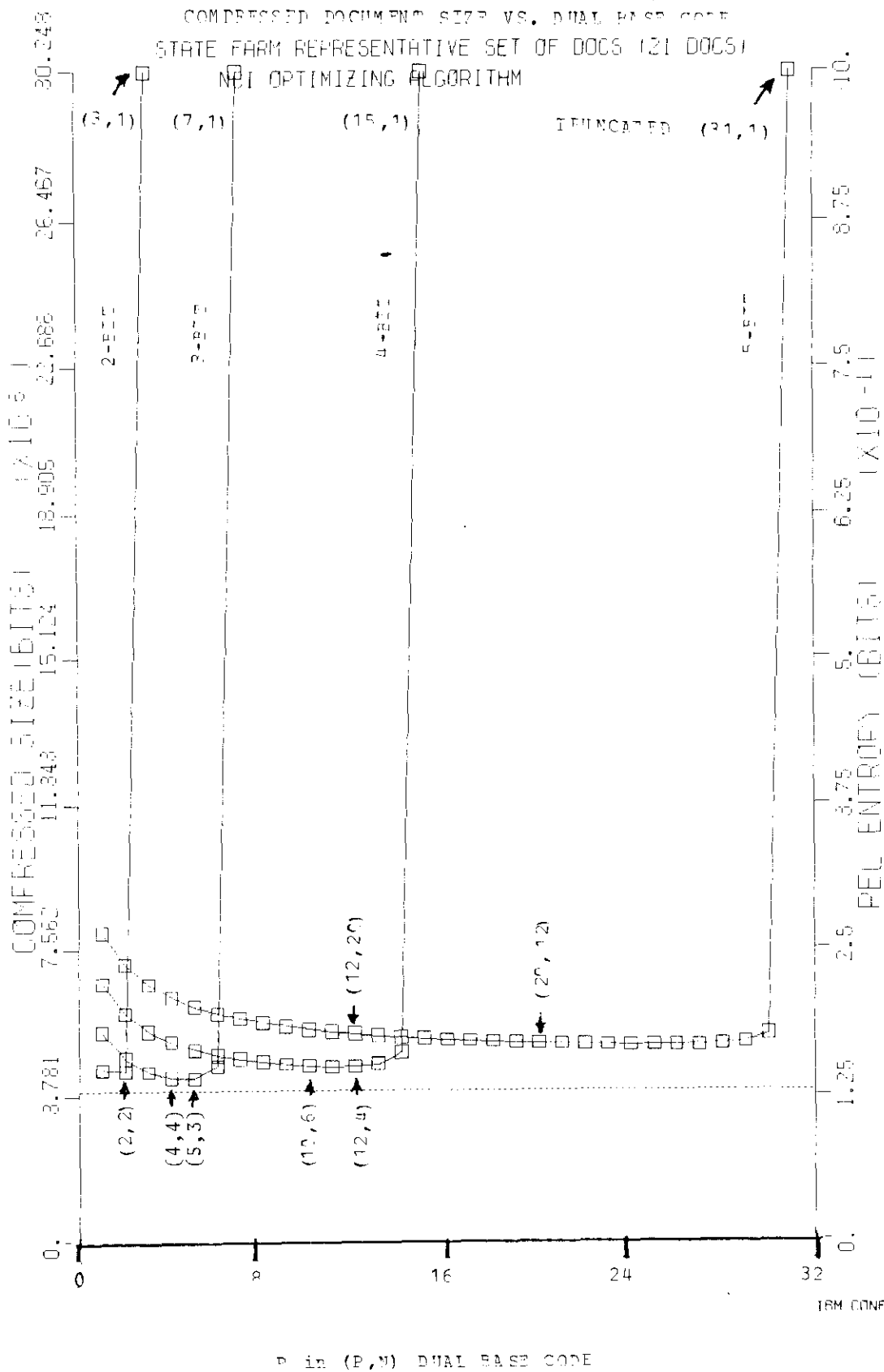
When using a variable word length, there must be a way of separating the individual code words out of a data stream. This is accomplished by restricting the low order digit (4 bits) of each code word to take on only the values 0 - 9. The high order digits (4 bits each) are restricted to the hex values A-F. The highest order digit is prevented from taking on the hex value A (it is restricted to B-F) since a code word of a hex A is given a different interpretation. Normally, the code represents the length of a run of white pels (non-errors in the error image). It is assumed that a black pel (error) follows the run of white. Two black pels in succession are represented by a white (non-error) run of zero length.

No ordinary code word is permitted to begin with a hex A. The code word for the last run in a scan line that does not have a black pel (error) at its end uses the hex A. This code of A indicates that the rest of the current scan line is all zeros (non-errors). Note that a blank line generally produces a line of non-errors and is indicated by a code of A. The total code, along with the run-lengths that can be encoded by the code words of various lengths are shown in Figure 13.

The compressor can be used with less than full pages (fields of images). A code word consisting of 4 hex F's is used to indicate the end of the partial image field. If padding is required to end on an integral byte boundary, five F's are used.

Decompression starts by separating out each code word. The encoded run length is used to construct a run of the proper length of all white followed by a black pel. This is repeated scan line by scan line to reconstruct the error image. The error image is then used to correct outputs from the predictor.

Starting with the all white scan line assumed to precede the first scan line, the first pel of the page is of course predicted to be white, which is compared with the first pel of the error image to reconstruct the first pel of the original image. This pel is then stored to be used for later predictions. The above



18M CONFIDENTIAL

Fig. 12. Compressed Document Size vs. Dual Base Code for Representative Set with NCI Predictor for all 2-bit, 3-bit, 4-bit, and 5-bit Dual Base Counters

process is repeated until the whole original image is reconstructed. Both the run length decoding and the prediction proceed concurrently.

APPENDIX A - NCI Algorithm

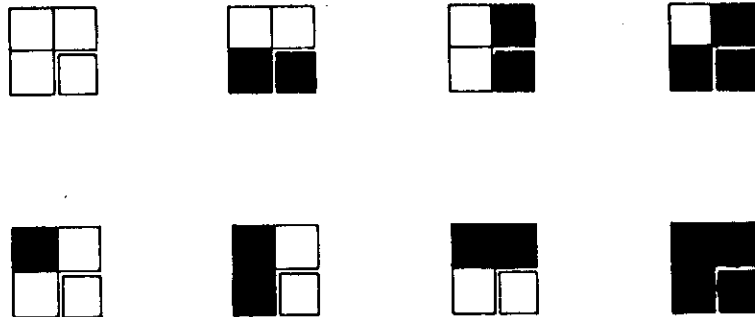
The NCI Algorithm used in the NCI prototype system at Los Gatos handles black-white images using the compression algorithm described as follows: The two step compression process proceeds as the image is first passed through a two-dimensional predictor. The output of the predictor is then passed to the run length encoder. Decompression is simply the inverse of this process.

The value of the current picture element (pel) is predicted based on the values of the immediately preceding pel in the current scan line, and the two adjacent pels in the previous scan line. Figure 13 shows the 8 predicted values for all possible combinations of the predicting pels. The output of the predictor is compared with the value of the current pel in the original image. If they are the same, a 0 bit is passed to the run length encoder. If they are different, an error has been made, and a 1 bit is passed to the run length encoder.

The result of this process can be thought of as a new image where the 1 bits represent black and the 0 bits represent white. This is called the error image. No compression has occurred up to this point; the error image has the same number of pels as were contained in the original image. The error image is just a transformation of the original image into a form that makes the run length encoder more efficient. One result of the prediction is that the error image has less black in it than did the original.

To start the prediction process, it is assumed that there is one scan line plus one pel which are all white preceding the first pel in the original image.

The run length coder is designed to be as compatible as possible with byte oriented IBM computers and still maintain high efficiency. It uses a variable word length, but the length can take on only the values 4, 8, 12, or 16 bits. At the end of a coded page, the



Predictor Algorithm

			0-9	0-9
	B-F	B-F	0-9	10-59
	A-F	A-F	0-9	60-359
B-F	A-F	A-F	0-9	360-2159

Run Length Code

Fig. 13. Predictor Table and Run Length Counter Coding for NCI Algorithm