

DOE/NV/1087a--T187

**UNLV
Information Science
Research Institute**

Quarterly Progress Report

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

T. A. Nartker
March 31, 1995

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED

MASTER

DISCLAIMER

**Portions of this document may be illegible
in electronic image products. Images are
produced from the best available original
document.**

Table of Contents

I. Board of Advisors Activity	1
II. Symposium Activity	1
III. Staff Activity	1
Recruiting	
Travel	
Papers accepted or presented	
IV. Document Analysis Program	2
OCR Devices	
OCR Test system	
OCR Databases/GT1	
OCR Databases/Foreign languages	
OCR Experiments	
OCR Technical reports/thesis	
Interaction with OCR vendors	
Interaction with OCR research organizations	
V. Text-Retrieval Program	3
TR Software systems	
TR Databases	
TR Experiments/Projects	
TR Technical reports/thesis	
Document Routing Project	
VI. Institute Activity	3
Institute visitors	
Institute seminars	
New agency contacts/research proposals	
VII. Goals Achieved This Quarter/Goals for Next Quarter	4
Appendix A. Final Program: SDAIR'95	
Appendix B. Status of Ground-truth Preparation Activities	

UNLV Information Science Research Institute

Quarterly Progress Report

T. A. Nartker
March 31, 1995

I. Board of Advisors Activity

A written report of both the Spring'94 and Fall'94 meetings is not yet available. Ms. Donna Harman, who acted as secretary for these meetings, has promised these reports soon.

The next Board meeting will be on Wednesday April 26th, after the SDAIR'95 Symposium.

II. Symposium Activity

The 1995 SDAIR Symposium is scheduled for Monday, Tuesday, & Wednesday, April 24-26 at the Desert Inn Hotel in Las Vegas. The four invited talks will be delivered by Dr. Richard Casey, Dr. Karen Sparck Jones, Dr. Abe Bookstein, and Dr. Hiromichi Fujisawa. Twenty five contributed papers and nine poster papers have been accepted. The Proceedings are being printed at this time. The Final Program is shown in Appendix A.

III. Staff Activity

Recruiting

None

Travel/ Meetings

In February, K. Taghva and J. Kanai traveled to San Jose, to present papers at SPIE'95. Kazem presented a paper titled "Evaluation of an Automatic Markup System" and Junichi presented a paper titled "Preliminary Evaluation of Histogram Based Binarization Algorithms."

In March, Shahram Latifi traveled to Snowbird Utah to present a paper "Compressed Image Manipulation," at the 1995 Data Compression Conference.

Papers accepted or presented

A paper by K. Taghva, J. Borsack, and A. Condit titled, "Effects of OCR Errors on Ranking and Feedback Using the Vector Space Model," was accepted for publication in Information Processing & Management.

A chapter for a Handbook of OCR and DIA titled "Information Retrieval and OCR" was submitted by K. Taghva, J. Borsack, and A. Condit.

IV. Document Analysis Program

OCR Devices

Testing for the ISRI fourth annual report of OCR accuracy has been completed. Eight organizations have submitted an OCR device(s) for this test. In addition to their binary engines, both Caere and HP have submitted greyscale engines. The international reputation of our annual tests is reflected in the number of devices submitted from foreign countries.

OCR Test system

Version #5 of our OCR experimental environment is complete and operational. It was used to conduct this years tests.

OCR Databases/GT1

All new datasets were completed and used for the fourth annual test of OCR accuracy. These include the DOE Sample#3 from the LSS prototype, the Magazine Sample, a Business Letter Sample, an English Newspaper Sample, and a Spanish Newspaper Sample.

Page 1 of Appendix B shows the ISRI methodology for preparing this ground-truth test data. Page 2 shows the status of all ISRI datasets on March 31. (The foreign language test data preparation shown is part of our DoD contract.)

OCR Databases/Foreign languages

As part of our contract with Fort Meade, we are preparing ground-truth text for a collection of newspaper articles printed in several foreign languages. We are collecting newspapers in Spanish, Chinese, and Japanese (as well as English).

As mentioned above, Appendix B gives an overview of our current ground-truth data preparation activities. Both the Spanish and the Chinese datasets are complete.

OCR Experiments

Testing, result analysis, and report writing are complete for our 1995 OCR Technology Assessment Report. This report is being printed at this time.

OCR Technical reports/thesis

Reports on "Prediction of OCR Accuracy" and on "Adaptive Restoration of Document Text" have been prepared and are included in our 1995 Annual Research Report.

Interaction with OCR vendors

None

Interaction with OCR research organizations

None

V. Text-Retrieval Program

TR Databases

We are continuing to collect relevance judgments for the queries associated with the residue (called GT1) of the LSS prototype database. We expect to complete all needed relevance judgments sometime this coming summer.

TR Experiments/Projects

Our text retrieval group has completed a prototype of the MANICURE OCR post-processing system. This system is available for testing by DOE (or TRW) staff. A report on MANICURE has been prepared and is being included in our 1995 Annual Research report.

TR Technical reports/thesis

"Autotag, A Tool for the Logical Markup of Documents" Allen Condit, M.S. Computer Science, May 1995. There are several other graduate students pursuing TR research projects. In our quarterly reports, we provide a summary of completed thesis projects but we do not provide interim progress reports for these projects.

VI. Institute Activity

LSS working group meetings & report

During the Fall quarter, T. Nartker and K. Taghva participated with the LSS Technical Working group.

Institute visitors

Our most important visits during the Winter quarter were representatives from Argonne National Labs. They expressed interest in using Manicure to recover the text from a DOE database of documents involving human radiation experiments. They also expressed interest in building a team, including ISRI staff, to build the LSS.

Institute visitors this quarter:

<u>Date</u>	<u>Visitor</u>	<u>Agency</u>
2/09/95	Bob Kero & Lucian Russell	Argonne National Labs
2/14/95	Representatives from ARMA	Las Vegas Chapter, Am. Records Mgt. Assn.

Institute seminars

None

New agency contacts/ new research proposals

Our foreign language work for Fort Meade has been suspended. We are awaiting continued funding for this project.

VII. Goals Achieved/Goals for Next Quarter

Goals from last quarter:

- 1) The final program for SDAIR'95 is shown in Appendix A. The Proceedings have been sent to the printer.
- 2) No new Affiliate program members were added during the Winter quarter.
- 3) Our 1995 Annual Report is complete and has gone to the printer.
- 4) All preparation of data samples is complete. (see Appendix C)
We expect to complete archival of these new datasets during the Spring quarter.
- 5) A version of MANICURE is complete and ready for initial testing.
- 6) Work has continued in collecting relevance judgments.

Goals for next quarter:

- 1) Select committee members and invited speakers for SDAIR'96.
- 2) Continue to recruit new members for the affiliates program.
- 3) Select datasets for our 1996 Annual OCR Test.
- 4) Continue work with the LSS Technical Working Group.
- 5) Investigate the use of MANICURE in the RDMS.
- 6) Continue work on relevance judgments for queries associated with the documents in GT1.

APPENDIX A.

SDAIR'95: Final Program

Appendix A Removed.
cont. Brochure
na

APPENDIX B.

ISRI Methodology for Preparing Ground-Truth Test Data and
Status of Ground-Truth Data Preparation Activities

ISRI METHODOLOGY FOR PREPARING GROUND-TRUTH OCR DATA

- 1. Acquire document sample**
 - Choose document class
 - Identify source for documents
 - Choose selection strategy and method of acquisition
 - Obtain documents

- 2. Select page sample**
 - Choose page (or partial page) sampling strategy
 - Sample and prepare pages
 - Choose zone types

- 3. Train data entry personnel**
 - Select & acquire tools
 - Train data entry staff

- 4. Scan pages**
 - Determine scanning variables (i.e. threshold & page placement)
 - Scan all pages (usually at 200, 300, 400, &GS)
 - Verify images
 - Quality control

- 5. Archive document & page collections**

- 6. Manually zone images**

- 7. Prepare Truth text**
 - Prepare multiple manual entries
 - Prepare text from ISRI voting algorithm
 - Resolve multiple manual entry & voting differences

- 8. Archive images,
zone information,
all manual & voting truth input, &
ground-truth text**

ISRI Ground-Truth Databases (As of 3/31/95)