Proceedings

# Fifth Annual Symposium
# on Document Analysis
# and Information Retrieval

April 15 - 17, 1996

Las Vegas, Nevada

# UNLV
UNIVERSITY OF NEVADA LAS VEGAS

Sponsored by:

## Information Science Research Institute

and

## Howard R. Hughes College of Engineering

University of Nevada, Las Vegas
4505 Maryland Parkway, Box 454021
Las Vegas, Nevada 89154-4021

# Table of Contents

## INVITED PAPERS

## CONTRIBUTED PAPERS