# TIPSTER TEXT PROGRAM

# PHASE II

## Proceedings of a Workshop held at

## Vienna, Virginia

## May 6-8, 1996

## Sponsored by:

## Defense Advanced Research Projects Agency

This document contains copies of reports prepared for the DARPA TIPSTER Text
Program - Phase II Workshop. Included are reports from DARPA sponsored program and
other materials prepared for use at the workshop.

**APPROVED FOR PUBLIC RELEASE
DISTRIBUTION UNLIMITED**

# Table of Contents

# Table of Contents

## Section C - Projects Employing TIPSTER Technologies

# Table of Contents

## Section D - Research in Information Extraction & Document Detection

# Table of Contents

## Section E - TIPSTER Architecture

# Table of Contents

## Section F - Evaluating the Technologies:
## The Text REtrieval Conferences (TREC)

## Section G - Evaluating the Technologies:
## The Message Understanding Conferences (MUC)

## Section H - Multi-Lingual Entity Task

# Table of Contents

## Section I - TIPSTER References