# HUMAN LANGUAGE TECHNOLOGY

Proceedings of a Workshop held at
Plainsboro, New Jersey

March 8-11, 1994

Sponsored by:
Advanced Research Projects Agency
Software & Intelligent Systems Technology Office

# TABLE OF CONTENTS

# AUTHOR INDEX