# Automatic Indexing:
# A State-of-the-Art Report

Mary Elizabeth Stevens

Center for Computer Sciences and Technology
National Bureau of Standards
Washington, D. C. 20234

Foreword
(1970 Edition)

Widespread interest in the use of computers in automatic indexing created a demand
for this publication that led to a recent exhaustion of all stock. While updating and revision
would have been desirable, other demands have prescribed reissuance with additional mate-
rial added as appendices. These are a paper updating the field through September 1966
(Appendix B), and bibliographic citations, pertinent to the subjects in the original text,
through August 1969 (Appendix C).

Lewis M. Branscomb
Director

The Research Information Center and Advisory Service on Information Processing, (RICASIP), which is jointly supported by the National Science Foundation and the National Bureau of Standards, is engaged in a continuing program to collect information and maintain current awareness about research and development activities in the field of information processing and retrieval. An important responsibility of RICASIP is the preparation of state-of-the-art reviews on topics of current interest in various areas of this broad field.

This report is one of a series intended as contributions toward improved interchange of information among those engaged in research and development in this field. The report considers new uses of machines and automatic data processing procedures for the compilation and generation of indexes to the scientific and technical literature.

A. V. Astin, Director

# Contents