Manfred R. Schroeder

# Computer Speech

## Recognition, Compression, Synthesis

Second Edition

With Introductions to Hearing and Signal Analysis
and a Glossary of Speech and Computer Terms

With 90 Figures

Springer

# Contents