# Text-to-Speech Synthesis

PAUL TAYLOR
*University of Cambridge*

# Contents