# *Inferring Phylogenies*

## Joseph Felsenstein
*University of Washington*

**Sinauer Associates, Inc.** • Publishers
Sunderland, Massachusetts

# Contents