# Biological Sequence Analysis

# Lecture Notes from CSE 527

Martin Tompa
Department of Computer Science and Engineering
University of Washington
Box 352350
Seattle, Washington, U.S.A. 98195-2350

Winter 2000

# Contents