# BioInformatics

## A COMPUTING PERSPECTIVE

Shuba Gopal

Anne Haake

Rhys Price Jones

Paul Tymann

# Contents