# Computational Molecular Evolution

Ziheng Yang ([z.yang@ucl.ac.uk](z.yang@ucl.ac.uk))
University College London

**Preface**

# Part I: Modeling Molecular Evolution

# Part II: Phylogeny Reconstruction

# Part III: Advanced Topics

## CHAPTER 9 Simulating Molecular Evolution

## CHAPTER 10 Perspectives

## Appendixes

## Reference
## Index