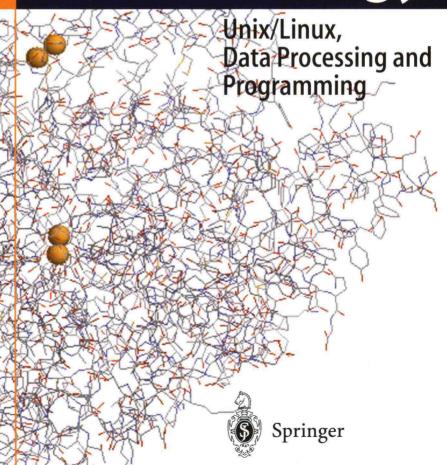


Computational Biology



Computational Biology —

Unix/Linux,
Data Processing and Programming

With 19 Figures and 12 Tables



Dr. Röbbe Wünschiers University of Cologne Institute for Genetics Weyertal 121 50931 Köln Germany

ISBN 3-540-21142-X Springer-Verlag Berlin Heidelberg New York

Library of Congress Control Number: 2004102411

This work is subject to copyright. All rights reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

Springer-Verlag is a part of Springer Science + Business Media springeronline.com

Springer-Verlag Berlin Heidelberg 2004 Printed in Germany

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Cover design: Design & Production, Heidelberg Typesetting: Camera ready by the author 31/3150WI - 5 4 3 2 1 0 - Printed on acid-free paper

Röbbe Wünschiers

Computational Biology -Unix/Linux, Data Processing and Programming

Springer

Berlin
Heidelberg
New York
Hong Kong
London
Milan
Paris
Tokyo

此为试读,需要完整PDF请访问: www.ertongbook.com

Dedicated to Károly Nagy and to the Open-Source Community

Foreword

A shift in culture

Only a decade ago, the first thing a molecular biologist would have had to learn when he or she started the lab work was how to handle pipettes, extract DNA, use enzymes and clone a gene. Now, the first thing that he or she should learn is how to handle databases and to extract all the information that is already known about the gene that he or she wants to study. In all likelihood, he or she will find that the gene has already been sequenced from several organisms, that it was recovered in a variety of EST projects, that expression data are available from microarray and SAGE studies, that it was included in linkage studies, that proteomics data are rapidly accumulating, that lists of interacting proteins are being compiled, that domain structure data are available and that it is part of a network of genetic interactions which is intensively modelled. He or she will discover that all this information resides in many different databases with different data formats and with different levels of analyses and linking. Starting to work on this gene will make sense only, if all this information is put together in a project-specific manner and set into the context of what is known about related genes and processes. At this point he or she may decide to walk up to the bioinformatics group in house and ask for help with arranging the data in a useful manner. This will then turn into the first major frustration in his or her career, since the last thing a scientific bioinformatics group wants to do is to provide a service for data retrieval and management.

Molecular biology is currently going through a dramatic cultural shift. The daily business of pipetting and gel running has increasingly to be complemented with data compiling and processing. Large lists of data are produced by sequencers, microarray experiments or real-time PCR machines every day. Working with data lists has become as important as extracting DNA. Every bench scientist needs proficiency in computing; discoveries are made both at the bench and on the screen. It would be completely wrong to think that computers in molecular biology are the business of bioinformaticians only.

VIII Foreword

Bioinformatics has become a scientific discipline of its own and should not be considered to be a service provider. The day-to-day computing will always have to be done by the experimentalist himself or herself.

Of course, there are now also a lot of helpful and fancy program packages for the bench scientist; but these will only perform routine tasks and all too often they are only poorly compatible. A scientist needs the freedom to develop his or her own ideas and to link things that have previously not been linked. Being able to go back to the basics of computing and programming is therefore a vital skill for the experimentalist, as important as making buffers and setting up enzyme reactions. It allows him or her to handle and analyze the data in exactly the way it is required for the project and to pursue new avenues of research, rather than trotting old paths.

Unix is the key to basic computing. If one is used to Windows or Mac operating systems, this might at first sound like going back into the stone age; but the dramatic recent shift of at least the Mac operating system to a Unix base should teach us otherwise. Unix is here to stay and it allows the largest flexibility for bioinformatics applications. Those who have learned Unix will soon discover the myriad of little "progies" that are available from colleagues all over the world and that can make life much easier in the lab.

This book gives exactly the sort of introduction into Unix, Unix-based operating systems and programming languages that will be a key competence for experimentally working molecular biologists and that will make all the difference for the successful projects of the future. It has been written by a bench scientist, specifically with the needs of molecular biologists in mind. It can be used either for self-teaching or in practical courses. Every group leader should hand over this book to new students in the lab, together with their first set of pipettes.

Cologne/Germany, January 2004 Prof. Dr. Diethard Tautz (Institut für Genetik der Universität zu Köln)

Preface

Welcome on board!

With this book I would like to invite you, the scientist, to a journey through terminals and program codes. You are welcome to put aside your pipette, culture flask or rubber boots for a while, make yourself comfortable in front of a computer (do not forget your favourite hot alcohol-free drink) and learn some unixing and programming. Why? Because we are living in the information age and there is a huge amount of biological knowledge and databases out there. They contain information on almost everything: genes and genomes, rRNAs, enzymes, protein structures, DNA-microarray experiments, single organisms, ecological data, the tree of life and endless more. Furthermore, nowadays many research apparatuses are connected to computers. Thus, you have electronic access to your data. However, in order to cope with all this information you need some tools. This book will provide you with the skills to use these tools and to develop your own tools, i.e. it will introduce Unix and its derivatives (Linux, Mac OS X, CygWin, etc.) and programming (shell programming, awk, perl). These tools will make you independent of the way in which other people make you process your data - in the form of application software. What you want is open functionality. You want to decide how to process (e.g. analyze, format, save, correlate) data and you want it now - not waiting for the lab programmer to treat your request; and you know it best - you understand your data and your demands. This is what open functionality stands for, and both Linux and programming languages can provide it to you.

I started programming on a Casio PB-100 hand-held built in 1983. It can store 10 small Basic programs. The accompanying book was entitled "Learn as you go" and, indeed, in my opinion this is the best way to learn programming. My first contact to Unix was triggered by the need to copy data files from a Unix-driven Bruker EPR-Spectrometer onto a floppy disk. The real challenge started when I tried to import the files to a data-plotting program on the PC. While the first problem could be solved by finding the right page in a Unix manual, the latter required programming skills – Q-Basic at that time. This

problem was minor compared to the trouble one encounters today. A common problem is to feed one program with the output of another program: you might have to change lines to columns, commas to dots, tabulators to semicolons, uppercase to lowercase, DNA to RNA, FASTA to GenBank format and so forth. Then there is that huge amount of information out there in the web, which you might need to bring into shape for your own analysis.

You and This Book – This book is written for the total beginner. You need not even to know what a computer is, though you should have access to one and find the power switch. The book is the result of a) the way how I learned to work with Unix, its derivatives and its numerous tools and b) a lecture which I started at the Institute for Genetics at the University of Cologne/Germany. Most programming examples are taken from biology; however, you need not be a biologist. Except for two or three examples, no biological knowledge is necessary. I have tried to illustrate almost everything practically with so-called terminals and examples. You should run these examples. Each chapter closes with some exercises. Brief solutions can be found at the end of the book.

Why Linux? – This book is not limited to Linux! All examples are valid for Unix or any Unix derivative like Mac OS X, Knoppix or the free Windows-based CygWin package, too. I chose Linux because it is open source software: you need not invest money except for the book itself. Furthermore, Linux provides all the great tools Unix provides. With Linux (as with all other Unix derivatives) you are close to your data. Via the command line you have immediate access to your files and can use either publicly available or your own designed tools to process these. With the aid of pipes you can construct your own data-processing pipeline. It is great.

Why awk and perl? — awk is a great language for both learning programming and treating large text-based data files (contrary to binary files). To 99% you will work with text-based files, be it data tables, genomes or species lists. Apart from being simple to learn and having a clear syntax, awk provides you with the possibility to construct your own commands. Thus, the language can grow with you as you grow with the language. I know bioinformatic professionals entirely focusing on awk. perl is much more powerful but also more unclear in its syntax (or flexible, to put it positively), but, since awk was one basis for developing perl, it is only a small step to go once you have learned awk — but a giant leap for your possibilities. You should take this step. By the way, both awk and perl run on all common operating systems.

Acknowledgements – Special thanks to Kristina Auerswald, Till Bayer, Benedikt Bosbach and Chris Voolstra for proofreading, and all the other students for encouraging me to bring these lines together.

Hürth/Germany, January 2004

Contents

| Pa | rt I | Whetting Your Appetite |
|----|---|---|
| 1 | Int. 1.1 1.2 1.3 1.4 1.5 1.6 1.7 | roduction 3 Information 3 What Is Linux? 3 What Is Shell Programming? 4 What Is sed? 4 What Is awk? 4 What Is per1? 5 Prerequisites 5 Conventions 6 |
| Pa | rt II | Computer and Operating Systems |
| 2 | 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 | ix/Linux 9 What Is a Computer? 9 Some History 10 2.2.1 Versions of Unix 11 2.2.2 The Rise of Linux 12 2.2.3 Why a Penguin? 14 2.2.4 Linux Distributions 14 X-Windows 14 2.3.1 How Does It Work? 15 Unix/Linux Architecture 16 What Is the Difference Between Unix and Linux? 17 What Is the Difference Between Unix/Linux and Windows? 17 What Is the Difference Between Unix/Linux and Mac OS X? 18 One Computer, Two Operating Systems 18 2.8.1 VMware 19 2.8.2 CygWin 19 |

| XII | Contents |
|-----|----------|
| | |

| | 2.9 2.10 | 2.8.4 Others | 19 20 20 21 21 22 |
|----|--|--|----------------------------------|
| Pa | rt III | Working with Unix/Linux | |
| 3 | The | Tilst Todell | 27 |
| | 3.1 | 3.1.1 Working with CygWin or Mac OS X | 27 27 28 28 |
| | 3.2 | 0.1.0 Working on a recine county and | 30 31 32 32 33 |
| | 3.3 Exerc | Logout | 33 34 |
| 4 | 4.1 4.2 4.3 4.4 4.5 4.6 4.7 Exerc | Browsing Files and Directories File Attributes Special Files: . and Protecting Files and Directories 4.4.1 Directories 4.4.2 Files 4.4.3 Examples 4.4.4 Changing File Attributes 4.4.5 Extended File Attributes File Archives File Compression Searching for Files cises | 49 50 |
| 5 | 5.1 5.2 5.3 5.4 | Downloading the Programs via FTP 5.1.1 Downloading BLAST 5.1.2 Downloading ClustalW Installing BLAST Running BLAST Installing ClustalW | 53 54 55 57 57 59 |

| | | | Contents | XIII |
|---|--------|---|--------------------------------|------|
| | 5.5 | Running ClustalW | | F0 |
| | | cises | | |
| | Direct | 0.000 | | . 00 |
| 6 | Wor | king with Text | 5 | . 63 |
| | 6.1 | A Quick Start: cat | | . 64 |
| | | 6.1.1 Text Sorting | | |
| | | 6.1.2 View File Beginning or End | | |
| | | 6.1.3 Scrolling Through Files | | . 67 |
| | | 6.1.4 Character, Word and Line Counting | | |
| | | 6.1.5 Finding Text | | |
| | 0.0 | 6.1.6 Text File Comparisons | | |
| | 6.2 | pico | | |
| | 6.3 | vi and vim | | |
| | | 6.3.1 Immediate Takeoff | | |
| | | 6.3.2 Starting vi | | |
| | | | | |
| | | G ==================================== | | |
| | | 8 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 | | |
| | | 6.3.6 Save and Quit | | |
| | | 6.3.8 Search and Replace | | |
| | Exer | cises | | |
| | | | | . 10 |
| 7 | Usin | g the Shell | | |
| | 7.1 | What Is the Shell? | Leta e e e e e e e e e e e e e | . 81 |
| | 7.2 | Different Shells | | . 82 |
| | 7.3 | Setting the Default Shell | | |
| | 7.4 | Useful Shortcuts | | |
| | 7.5 | Redirections | | |
| | 7.6 | Pipes | | |
| | 7.7 | Lists of Commands | | |
| | 7.8 | Aliases | | |
| | 7.9 | Scheduling Commands | *********** | . 89 |
| | 7.10 | Wildcards | | |
| | 7.11 | Processes | | |
| | | 7.11.1 Checking Processes | | |
| | | | | |
| | Ever | 7.11.3 Killing Processes | | |
| | DVCL | 2000 | | . 96 |
| 8 | Shell | Programming | | . 97 |
| | 8.1 | Script Files | | |
| | 8.2 | Modifying the Path | | |
| | 8.3 | Variables | | |
| | 8.4 | Input and Output | | 102 |

| | | 8.4.1 | echo 103 |
|---|-------|----------|--|
| | | 8.4.2 | Here Documents: << |
| | | 8.4.3 | read and line104 |
| | | 8.4.4 | Script Parameters |
| | 8.5 | | ntions and Expansions |
| | 0.0 | 8.5.1 | Variable Substitution |
| | | 8.5.2 | Command Expansion |
| | 8.6 | | g |
| | 0.0 | 8.6.1 | Escape Character |
| | | 8.6.2 | Single Quotes |
| | | 8.6.3 | Double Quotes |
| | 8.7 | | ns – Flow Control |
| | 0 | 8.7.1 | ifthenelifelsefi |
| | | 8.7.2 | test111 |
| | | 8.7.3 | whiledodone |
| | | 8.7.4 | untildodone |
| | | 8.7.5 | forindodone |
| | | 8.7.6 | caseinesac |
| | | 8.7.7 | selectindo |
| | 8.8 | Debugg | ing119 |
| | | 8.8.1 | bash -xv120 |
| | | 8.8.2 | trap121 |
| | 8.9 | Exampl | les |
| | | 8.9.1 | Check for DNA as File Content122 |
| | | 8.9.2 | Time Signal |
| | | 8.9.3 | Select Files to Archive |
| | | 8.9.4 | Remove Spaces |
| | Exerc | eises | $\dots \dots $ |
| | | | |
| 9 | | ılar Exp | pressions |
| | 9.1 | Using F | Regular Expressions |
| | 9.2 | | Pattern and Examples |
| | | 9.2.1 | Single-Character Meta Characters |
| | | 9.2.2 | Quantifiers |
| | | 9.2.3 | Grouping |
| | | 9.2.4 | Anchors |
| | | 9.2.5 | Escape Sequences |
| | | 9.2.6 | Alternation |
| | | 9.2.7 | Back References |
| | | 9.2.8 | Character Classes |
| | | 9.2.9 | Priorities |
| | | 9.2.10 | egrep Options |
| | 9.3 | 0 | r Expressions and vim |
| | Exerc | cises | $\dots \dots $ |

| | | | Contents | XV |
|----|-------------|---|---------------------------------------|-------|
| 10 | Sod | | | H X H |
| 10 | 10 Sed | | | |
| | 10.1 | Getting Started | | |
| | 10.2 10.3 | How sed Works | | |
| | 10.0 | 10.3.1 Pattern Space | | |
| | | 10.3.2 Hold Space | | |
| | 10.4 | sed Syntax | | |
| | 10.1 | 10.4.1 Addresses | | |
| | | 10.4.2 sed and Regular Expressions | | |
| | 10.5 | Commands | | |
| | | 10.5.1 Substitutions | | |
| | | 10.5.2 Transliterations | | |
| | | 10.5.3 Deletions | | |
| | | 10.5.4 Insertions and Changes | | |
| | | 10.5.5 sed Script Files | | |
| | | 10.5.6 Printing | | |
| | | 10.5.7 Reading and Writing Files | | |
| | | 10.5.8 Advanced sed | | 157 |
| | 10.6 | Examples | | 157 |
| | | 10.6.1 Gene Tree | | 157 |
| | | 10.6.2 File Tree | | |
| | | 10.6.3 Reversing Line Order | | |
| | Exerc | ises | | 160 |
| | | | | |
| Pa | rt IV | Programming | | |
| - | | | | |
| 11 | Awk | ******************************* | 90 9 190 91 91 91 92 9 93 8 196 9 | 163 |
| | 11.1 | Getting Started | ******* | 164 |
| | 11.2 | awk's Syntax | | 165 |
| | 11.3 | Example File | | 166 |
| | 11.4 | Patterns | ********* | 167 |
| | | 11.4.1 Regular Expressions | | 167 |
| | | 11.4.2 Pattern-Matching Expressions | | |
| | | 11.4.3 Relational Character Expressions | | |
| | | 11.4.4 Relational Number Expressions | | |
| | | 11.4.5 Mixing and Conversion of Numbers and C | | |
| | | 11.4.6 Ranges | | |
| | 3312 | 11.4.7 BEGIN and END | | |
| | 11.5 | Variables | | |
| | | 11.5.1 Assignment Operators | | |
| | | 11.5.2 Increment and Decrement | | |
| | | 11.5.3 Predefined Variables | | |
| | | 11.5.4 Arrays | **** | 181 |
| | | 11.5.5 Shell Versus awk Variables | * * * * * * * * * * * * * * * * * * * | 185 |

| | 11.6 | Scripts and Executables | | | |
|----|-----------------------|---|--|--|--|
| | 11.7 | | ns – Flow Control | | |
| | | 11.7.1 | ifelse | | |
| | | 11.7.2 | while | | |
| | | 11.7.3 | dowhile | | |
| | | 11.7.4 | for | | |
| | | 11.7.5 | Leaving Loops | | |
| | 11.8 | Actions | 5 | | |
| | | 11.8.1 | Printing | | |
| | | 11.8.2 | Numerical Calculations | | |
| | | 11.8.3 | String Manipulation | | |
| | | 11.8.4 | System Commands | | |
| | | 11.8.5 | User-Defined Functions | | |
| | 11.9 | Input with getline | | | |
| | | 11.9.1 Reading from a File | | | |
| | | 11.9.2 | Reading from the Shell | | |
| | 11.10 | Examp | les | | |
| | | 11.10.1 | Sort an Array by Indices | | |
| | | | Sum up Atom Positions | | |
| | | | Translate DNA to Protein | | |
| | | | Calculate Atomic Composition of Proteins | | |
| | | 11.10.5 | Calculate Distance Between Cysteines | | |
| | Exercises | | | | |
| | Exerc | ises | | | |
| 10 | | | | | |
| 12 | Perl . | | | | |
| 12 | Perl . | Intentio | | | |
| 12 | Perl . 12.1 12.2 | Intentic | 221 on of this Chapter | | |
| 12 | Perl . | Intentic Runnin Variable | | | |
| 12 | Perl . 12.1 12.2 | Intentic Runnin Variable 12.3.1 | | | |
| 12 | Perl . 12.1 12.2 | Intention Runnin Variable 12.3.1 12.3.2 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 | | |
| 12 | Perl . 12.1 12.2 | Intention Runnin Variable 12.3.1 12.3.2 12.3.3 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intention Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 | | |
| 12 | Perl . 12.1 12.2 | Intention Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 ifelseif .else 234 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 if . elseif . else 234 unless, die and warn 234 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intention Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 if . elseif . else 234 unless, die and warn 234 while 235 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 12.4.4 | on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 ifelseifelse 234 unless, die and warn 234 while 235 dowhile 235 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intention Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 if. elseif. else 234 unless, die and warn 234 while 235 do. while 235 until 235 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 12.4.4 12.4.5 12.4.6 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 if . elseif . else 234 unless, die and warn 234 while 235 do . while 235 until 235 do . until 235 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentice Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 12.4.4 12.4.5 12.4.6 12.4.7 | 221 on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 ifelseifelse 234 unless, die and warn 234 while 235 do. while 235 until 235 for 236 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 12.4.4 12.4.5 12.4.6 12.4.7 12.4.8 | on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 ifelseif.else 234 unless, die and warn 234 while 235 dowhile 235 dountil 235 for 236 foreach 236 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 12.4.4 12.4.5 12.4.6 12.4.7 12.4.8 12.4.9 | on of this Chapter 221 on of this Chapter 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 in on Flow Control 234 if olseif else 234 unless, die and warn 234 while 235 do. while 235 do. until 235 for 236 foreach 236 Controlling Loops 237 | | |
| 12 | Perl . 12.1 12.2 12.3 | Intentic Runnin Variable 12.3.1 12.3.2 12.3.3 12.3.4 Decision 12.4.1 12.4.2 12.4.3 12.4.4 12.4.5 12.4.6 12.4.7 12.4.8 12.4.9 | on of this Chapter 221 g perl Scripts 222 es 222 Scalars 223 Arrays 224 Hashes 229 Built-in Variables 233 ns - Flow Control 234 ifelseif.else 234 unless, die and warn 234 while 235 dowhile 235 dountil 235 for 236 foreach 236 | | |

| | | Contents A | V 11 |
|------------|--------|--|------|
| | | 12.5.3 Files | 240 |
| | 12.6 | Data Output | |
| | 12.0 | 12.6.1 print, printf and sprintf | |
| | | 12.6.2 Here Documents: << | |
| | | 12.6.3 Files | |
| | 12.7 | Hash Databases | |
| | 12.8 | Regular Expressions | |
| | 12.0 | 12.8.1 Special Escape Sequences | |
| | | 12.8.2 Matching: m// | |
| | 12.9 | String Manipulations | |
| | 12.9 | | |
| | | | |
| | | | |
| | 10.10 | | |
| | | Calculations | |
| | | Subroutines | |
| | | Packages and Modules | |
| | | Bioperl | |
| | | You Want More? | |
| | 12.15 | Examples | |
| | | 12.15.1 Reverse Complement DNA | |
| | | 12.15.2 Calculate GC Content | |
| | | 12.15.3 Restriction Enzyme Digestion | |
| | - | 12.15.4 Levenshtein Distance of Sequences | |
| | Exerc | ises | 263 |
| A | Anno | endix | 265 |
| A . | A.1 | Keyboard Shortcuts | |
| | A.2 | Boot Problems | |
| | A.3 | Optimizing the Bash Shell | |
| | A.5 | A.3.1 Startup Files | |
| | | A.3.2 A Helpful Shell Prompt | |
| | A.4 | Text File Conversion (Unix \leftrightarrow DOS) | |
| | A.5 | Devices | |
| | A.6 | Mounting Filesystems | |
| | A.0 | | |
| | | the state of the s | |
| | | | |
| | A.7 | The second secon | |
| | | Nucleotide and Protein Codes | |
| | A.8 | Special Characters | 272 |
| Sol | utions | S | 273 |
| | autom | × | -,0 |
| Re | ferenc | es | 279 |
| | | | 2.2 |
| Inc | lex | | 281 |