

*Understand Your Data and
Be More Productive*

3rd Edition
For Perl, PHP, Java,
.NET, Ruby, and More!



Mastering

Regular Expressions

O'REILLY®

Jeffrey E.F. Friedl

*Understand Your Data and
Be More Productive*

3rd Edition
For Perl, PHP, Java,
.NET, Ruby, and More!



Mastering

Regular Expressions

O'REILLY®

Jeffrey E.F. Friedl

Mastering Regular Expressions



This book will get you up to speed on the productivity secrets that will make your life easier: regular expressions. Well-crafted regular expressions can reduce hours of tedious labor to a 15-second solution. Now a standard feature in a wide range of languages and popular tools—including Perl, PHP, Java, Python, Ruby, MySQL, VB.NET, and C# (and any language using the .NET Framework)—regular expressions will allow you to code complex and subtle text processing that you never imagined could be automated.

Mastering Regular Expressions, Third Edition, now includes material on PHP and its powerful regular expressions. This edition has been updated throughout to reflect advances in other languages, including expanded in-depth coverage of Sun's `java.util.regex`, with special attention to the many differences between Java 1.4.2 and Java 1.5/1.6.

Topics include:

- A comparison of features among many languages and tools
- How a regular expression engine works
- Optimization (major savings available here!)
- Matching only what you want, not what you don't
- Sections and chapters on individual languages

Written in the lucid, entertaining tone that makes a complex topic become crystal clear to programmers, and filled with solutions to difficult real-world problems, *Mastering Regular Expressions*, Third Edition, offers a wealth of information that you can put to immediate use to craft elegant, time-saving solutions to a wide range of issues.

“If you use regular expressions as part of your professional work (even if you already have a good book on whatever language you’re programming in) I would strongly recommend this book to you.”

—Dr Chris Brown, Linux Format

“Mastering Regular Expressions is the definitive guide to the subject, and an outstanding resource that belongs on every programmer’s bookshelf. Ten out of Ten Horseshoes.”

—Jason Menard, Java Ranch

“There isn’t a better (or more useful) book available on regular expressions.”

—Zak Greant, Planet PHP

www.oreilly.com



Mastering Regular Expressions

Mastering Regular Expressions

Third Edition

Jeffrey E. F. Friedl

O'REILLY®

*Beijing • Cambridge • Farnham • Köln • Paris • Sebastopol • Taipei •
Tokyo*

Mastering Regular Expressions, Third Edition

by Jeffrey E. F. Friedl

Copyright © 2006, 2002, 1997 O'Reilly Media, Inc. All rights reserved.
Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly Media, Inc. books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (safari.oreilly.com). For more information contact our corporate/institutional sales department: 800-998-9938 or corporate@oreilly.com.

Editor: Andy Oram

Production Editor: Jeffrey E. F. Friedl

Cover Designer: Edie Freedman

Printing History:

January 1997:	First Edition.
July 2002:	Second Edition.
August 2006:	Third Edition.

Nutshell Handbook, the Nutshell Handbook logo, and the O'Reilly logo are registered trademarks of O'Reilly Media, Inc. *Mastering Regular Expressions*, the image of owls, and related trade dress are trademarks of O'Reilly Media, Inc. Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and O'Reilly Media, Inc. was aware of a trademark claim, the designations have been printed in caps or initial caps.

While every precaution has been taken in the preparation of this book, the publisher and author assume no responsibility for errors or omissions, or for damages resulting from the use of the information contained herein.



This book uses RepKover™, a durable and flexible lay-at binding.

ISBN: 0-596-52812-4
[M]

FOR Fumie
文枝

For putting up with me.
And for the years I worked on this book,
for putting up without me.

Table of Contents

Preface

1: Introduction to Regular Expressions

Solving Real Problems

Regular Expressions as a Language

 The Filename Analogy

 The Language Analogy

The Regular-Expression Frame of Mind

 If You Have Some Regular-Expression Experience

 Searching Text Files: Egrep

Egrep Metacharacters

 Start and End of the Line

 Character Classes

 Matching Any Character with Dot

 Alternation

 Ignoring Differences in Capitalization

 Word Boundaries

 In a Nutshell

 Optional Items

 Other Quantifiers: Repetition

 Parentheses and Backreferences

 The Great Escape

Expanding the Foundation

 Linguistic Diversification

 The Goal of a Regular Expression

 A Few More Examples

 Regular Expression Nomenclature

 Improving on the Status Quo

 Summary

Personal Glimpses

2: Extended Introductory Examples

About the Examples

A Short Introduction to Perl

Matching Text with Regular Expressions

Toward a More Real-World Example

Side Effects of a Successful Match

Intertwined Regular Expressions

Intermission

Modifying Text with Regular Expressions

Example: Form Letter

Example: Prettifying a Stock Price

Automated Editing

A Small Mail Utility

Adding Commas to a Number with Lookaround

Text-to-HTML Conversion

That Doubled-Word Thing

3: Overview of Regular Expression Features and Flavors

A Casual Stroll Across the Regex Landscape

The Origins of Regular Expressions

At a Glance

Care and Handling of Regular Expressions

Integrated Handling

Procedural and Object-Oriented Handling

A Search-and-Replace Example

Search and Replace in Other Languages

Care and Handling: Summary

Strings, Character Encodings, and Modes

Strings as Regular Expressions

Character-Encoding Issues

Unicode

Regex Modes and Match Modes

Common Metacharacters and Features

Character Representations

Character Classes and Class-Like Constructs

Anchors and Other “Zero-Width Assertions”
Comments and Mode Modifiers
Grouping, Capturing, Conditionals, and Control
Guide to the Advanced Chapters

4: The Mechanics of Expression Processing

Start Your Engines!

Two Kinds of Engines
New Standards
Regex Engine Types
From the Department of Redundancy Department
Testing the Engine Type

Match Basics

About the Examples
Rule 1: The Match That Begins Earliest Wins
Engine Pieces and Parts
Rule 2: The Standard Quantifiers Are Greedy

Regex-Directed Versus Text-Directed

NFA Engine: Regex-Directed
DFA Engine: Text-Directed
First Thoughts: NFA and DFA in Comparison

Backtracking

A Really Crummy Analogy
Two Important Points on Backtracking
Saved States
Backtracking and Greediness

More About Greediness and Backtracking

Problems of Greediness
Multi-Character “Quotes”
Using Lazy Quantifiers
Greediness and Laziness Always Favor a Match
The Essence of Greediness, Laziness, and Backtracking
Possessive Quantifiers and Atomic Grouping
Possessive Quantifiers, ?+, *+, ++, and {m,n}+
The Backtracking of Lookaround

- Is Alternation Greedy?
- Taking Advantage of Ordered Alternation
- NFA, DFA, and POSIX
 - “The Longest-Leftmost”
 - POSIX and the Longest-Leftmost Rule
 - Speed and Efficiency
 - Summary: NFA and DFA in Comparison
- Summary

5: *Practical Regex Techniques*

- Regex Balancing Act
- A Few Short Examples
 - Continuing with Continuation Lines
 - Matching an IP Address
 - Working with Filenames
 - Matching Balanced Sets of Parentheses
 - Watching Out for Unwanted Matches
 - Matching Delimited Text
 - Knowing Your Data and Making Assumptions
 - Stripping Leading and Trailing Whitespace
- HTML-Related Examples
 - Matching an HTML Tag
 - Matching an HTML Link
 - Examining an HTTP URL
 - Validating a Hostname
 - Plucking Out a URL in the Real World
- Extended Examples
 - Keeping in Sync with Your Data
 - Parsing CSV Files

6: *Crafting an Efficient Expression*

- A Sobering Example
 - A Simple Change—Placing Your Best Foot Forward
 - Efficiency Versus Correctness
 - Advancing Further—Localizing the Greediness

- Reality Check
- A Global View of Backtracking
 - More Work for a POSIX NFA
 - Work Required During a Non-Match
 - Being More Specific
 - Alternation Can Be Expensive
- Benchmarking
 - Know What You're Measuring
 - Benchmarking with PHP
 - Benchmarking with Java
 - Benchmarking with VB.NET
 - Benchmarking with Ruby
 - Benchmarking with Python
 - Benchmarking with Tcl
- Common Optimizations
 - No Free Lunch
 - Everyone's Lunch is Different
 - The Mechanics of Regex Application
 - Pre-Application Optimizations
 - Optimizations with the Transmission
 - Optimizations of the Regex Itself
- Techniques for Faster Expressions
 - Common Sense Techniques
 - Expose Literal Text
 - Expose Anchors
 - Lazy Versus Greedy: Be Specific
 - Split Into Multiple Regular Expressions
 - Mimic Initial-Character Discrimination
 - Use Atomic Grouping and Possessive Quantifiers
 - Lead the Engine to a Match
- Unrolling the Loop
 - Method 1: Building a Regex From Past Experiences
 - The Real "Unrolling-the-Loop" Pattern
 - Method 2: A Top-Down View

- Method 3: An Internet Hostname
- Observations
- Using Atomic Grouping and Possessive Quantifiers
- Short Unrolling Examples
- Unrolling C Comments
- The Freeflowing Regex
 - A Helping Hand to Guide the Match
 - A Well-Guided Regex is a Fast Regex
- Wrapup
- In Summary: Think!

7: Perl

- Regular Expressions as a Language Component
 - Perl's Greatest Strength
 - Perl's Greatest Weakness
- Perl's Regex Flavor
 - Regex Operands and Regex Literals
 - How Regex Literals Are Parsed
 - Regex Modifiers
- Regex-Related Perlisms
 - Expression Context
 - Dynamic Scope and Regex Match Effects
 - Special Variables Modified by a Match
- The `qr/.../` Operator and Regex Objects
 - Building and Using Regex Objects
 - Viewing Regex Objects
 - Using Regex Objects for Efficiency
- The Match Operator
 - Match's Regex Operand
 - Specifying the Match Target Operand
 - Different Uses of the Match Operator
 - Iterative Matching: Scalar Context, with `/g`
 - The Match Operator's Environmental Relations
- The Substitution Operator
 - The Replacement Operand

- The /e Modifier
- Context and Return Value
- The Split Operator
 - Basic Split
 - Returning Empty Elements
 - Split’s Special Regex Operands
 - Split’s Match Operand with Capturing Parentheses
- Fun with Perl Enhancements
 - Using a Dynamic Regex to Match Nested Pairs
 - Using the Embedded-Code Construct
 - Using `local` in an Embedded-Code Construct
 - A Warning About Embedded Code and `my` Variables
 - Matching Nested Constructs with Embedded Code
 - Overloading Regex Literals
 - Problems with Regex-Literal Overloading
 - Mimicking Named Capture
- Perl Efficiency Issues
 - “There’s More Than One Way to Do It”
 - Regex Compilation, the /o Modifier, `qr/.../`, and Efficiency
 - Understanding the “Pre-Match” Copy
 - The Study Function
 - Benchmarking
 - Regex Debugging Information
- Final Comments

8: Java

- Java’s Regex Flavor
 - Java Support for `\p{...}` and `\P{...}`
 - Unicode Line Terminators
- Using `java.util.regex`
- The `Pattern.compile()` Factory
 - Pattern’s `matcher` method
- The Matcher Object
 - Applying the Regex
 - Querying Match Results

- Simple Search and Replace
- Advanced Search and Replace
- In-Place Search and Replace
- The Matcher's Region
- Method Chaining
- Methods for Building a Scanner
- Other Matcher Methods
- Other Pattern Methods
 - Pattern's split Method, with One Argument
 - Pattern's split Method, with Two Arguments
- Additional Examples
 - Adding Width and Height Attributes to Image Tags
 - Validating HTML with Multiple Patterns Per Matcher
 - Parsing Comma-Separated Values (CSV) Text
- Java Version Differences
 - Differences Between 1.4.2 and 1.5.0
 - Differences Between 1.5.0 and 1.6

9: .NET

- .NET's Regex Flavor
 - Additional Comments on the Flavor
- Using .NET Regular Expressions
 - Regex Quickstart
 - Package Overview
 - Core Object Overview
- Core Object Details
 - Creating `Regex` Objects
 - Using `Regex` Objects
 - Using `Match` Objects
 - Using `Group` Objects
- Static "Convenience" Functions
 - Regex Caching
- Support Functions
- Advanced .NET
 - Regex Assemblies

Matching Nested Constructs

Capture Objects

10: PHP

PHP's Regex Flavor

The Preg Function Interface

“Pattern” Arguments

The Preg Functions

`preg_match`

`preg_match_all`

`preg_replace`

`preg_replace_callback`

`preg_split`

`preg_grep`

`preg_quote`

“Missing” Preg Functions

`preg_regex_to_pattern`

Syntax-Checking an Unknown Pattern Argument

Syntax-Checking an Unknown Regex

Recursive Expressions

Matching Text with Nested Parentheses

No Backtracking Into Recursion

Matching a Set of Nested Parentheses

PHP Efficiency Issues

The S Pattern Modifier: “Study”

Extended Examples

CSV Parsing with PHP

Checking Tagged Data for Proper Nesting

Index

Preface

This book is about a powerful tool called “regular expressions”. It teaches you how to use regular expressions to solve problems and get the most out of tools and languages that provide them. Most documentation that mentions regular expressions doesn’t even begin to hint at their power, but this book is about *mastering* regular expressions.

Regular expressions are available in many types of tools (editors, word processors, system tools, database engines, and such), but their power is most fully exposed when available as part of a programming language. Examples include Java and JScript, Visual Basic and VBScript, JavaScript and ECMAScript, C, C++, C#, elisp, Perl, Python, Tcl, Ruby, PHP, *sed*, and *awk*. In fact, regular expressions are the very heart of many programs written in some of these languages.

There’s a good reason that regular expressions are found in so many diverse languages and applications: they are extremely powerful. At a low level, a regular expression describes a chunk of text. You might use it to verify a user’s input, or perhaps to sift through large amounts of data. On a higher level, regular expressions allow you to master your data. Control it. Put it to work for you. To master regular expressions is to master your data.

The Need for This Book

I finished the first edition of this book in late 1996, and wrote it simply because there was a need. Good documentation on regular expressions just wasn’t available, so most of their power went untapped. Regular-expression documentation was available, but it centered on the “low-level view.” It seemed to me that they were analogous to showing someone the alphabet and expecting them to learn to speak.

The five and a half years between the first and second editions of this book saw the popular rise of the Internet, and, perhaps more than just coincidentally, a considerable expansion in the world of regular expressions. The regular expressions of almost every tool and language became more powerful and expressive. Perl, Python, Tcl, Java, and Visual Basic all got new regular-expression backends. New languages with regular expression support, like PHP, Ruby, and C#, were developed and became popular. During all this time, the basic core of the book — how to truly understand regular expressions and how to get the most from them — remained as important and relevant as ever.

Yet, the first edition gradually started to show its age. It needed updating to reflect the new languages and features, as well as the expanding role that regular expressions played in the Internet world. It was published in 2002, a year that saw the landmark

releases of `java.util.regex`, Microsoft's .NET Framework, and Perl 5.8. They were all covered fully in the second edition. My one regret with the second edition was that it didn't give more attention to PHP. In the four years since the second edition was published, PHP has only grown in importance, so it became imperative to correct that deficiency.

This third edition features enhanced PHP coverage in the early chapters, plus an all new, expansive chapter devoted entirely to PHP regular expressions and how to wield them effectively. Also new in this edition, the Java chapter has been rewritten and expanded considerably to reflect new features of Java 1.5 and Java 1.6.

Intended Audience

This book will interest anyone who has an opportunity to use regular expressions. If you don't yet understand the power that regular expressions can provide, you should benefit greatly as a whole new world is opened up to you. This book should expand your understanding, even if you consider yourself an accomplished regular-expression expert. After the first edition, it wasn't uncommon for me to receive an email that started "*I thought I knew regular expressions until I read *Mastering Regular Expressions*. Now I do.*"


Programmers working on text-related tasks, such as web programming, will find an absolute gold mine of detail, hints, tips, and *understanding* that can be put to immediate use. The detail and thoroughness is simply not found anywhere else.

Regular expressions are an idea—one that is implemented in various ways by various utilities (many, many more than are specifically presented in this book). If you master the general concept of regular expressions, it's a short step to mastering a particular implementation. This book concentrates on that idea, so most of the knowledge presented here transcends the utilities and languages used to present the examples.

How to Read This Book

This book is part tutorial, part reference manual, and part story, depending on when you use it. Readers familiar with regular expressions might feel that they can immediately begin using this book as a detailed reference, flipping directly to the section on their favorite utility. I would like to discourage that.

You'll get the most out of this book by reading the first six chapters as a story. I have found that certain habits and ways of thinking help in achieving a full understanding, but are best absorbed over pages, not merely memorized from a list.

The story that is the first six chapters form the basis for the last four, covering specifics of Perl, Java, .NET, and PHP. To help you get the most from each part, I've used cross references liberally, and I've worked hard to make the index as useful as possible. (Over 1,200 cross references are sprinkled throughout the book; they are often presented as "" followed by a page number.)