# WEBBOTS, SPIDERS, AND SCREEN SCRAPERS

## A GUIDE TO DEVELOPING *INTERNET AGENTS* WITH *PHP/CURL*

### MICHAEL SCHRENK

2ND EDITION

no starch press

# CONTENTS IN DETAIL

Webbots, Spiders, and Screen Scrapers, 2nd Edition
© 2012 Michael Schrenk

# 5
# ADVANCED PARSING WITH REGULAR EXPRESSIONS     49

# 6
# AUTOMATING FORM SUBMISSION     63

# PART IV: LARGER CONSIDERATIONS          263

## 26
## DESIGNING STEALTHY WEBBOTS AND SPIDERS          265

## 27
## PROXIES          273

## 31
## KEEPING WEBBOTS OUT OF TROUBLE 317

## A
## PHP/CURL REFERENCE 327

## B
## STATUS CODES 337

## C
## SMS GATEWAYS 341

## INDEX 345