DAVID MERTZ

# TEXT PROCESSING IN PYTHON

*Text Processing in Python* is an example-driven, hands-on tutorial that carefully teaches programmers how to accomplish numerous text processing tasks using the Python language. Filled with concrete examples, this book provides efficient and effective solutions to specific text processing problems and practical strategies for dealing with all types of text processing challenges.

*Text Processing in Python* begins with an introduction to text processing and contains a quick Python tutorial to get you up to speed. It then delves into essential text processing subject areas, including string operations, regular expressions, parsers and state machines, and Internet tools and techniques. Appendixes cover such important topics as data compression and Unicode. A comprehensive index and plentiful cross-referencing offer easy access to available information. In addition, exercises throughout the book provide readers with further opportunity to hone their skills either on their own or in the classroom. A companion Web site (http://gnosis.cx/TPiP) contains source code and examples from the book.

## ABOUT THE AUTHOR

DAVID MERTZ came to writing about programming via the unlikely route of first being a humanities professor. Along the way, he was a senior software developer, and now runs his own development company, Gnosis Software ("We know stuff!"). David writes regular columns and articles for IBM developerWorks, Intel Developer Network, O'Reilly ONLamp, and other publications.

Addison Wesley

FROM ADDISON-WESLEY

**DAVID MERTZ**

# TEXT PROCESSING IN PYTHON

Addison Wesley