

Lecture Notes in Computer Science

Edited by G. Goos and J. Hartmanis

96

James L. Peterson

Computer Programs for Spelling Correction:

An Experiment in Program Design



Springer-Verlag
Berlin Heidelberg New York 1980

Editorial Board

W. Brauer P. Brinch Hansen D. Gries C. Moler G. Seegmüller
J. Stoer N. Wirth

Author

James L. Peterson
The Department of Computer Sciences
The University of Texas
Austin, Texas 78712/USA

AMS Subject Classifications (1979): 68 A30
CR Subject Classifications (1974): 5.2, 3.42

ISBN 3-540-10259-0 Springer-Verlag Berlin Heidelberg New York
ISBN 0-387-10259-0 Springer-Verlag New York Heidelberg Berlin

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically those of translation, reprinting, re-use of illustrations, broadcasting, reproduction by photocopying machine or similar means, and storage in data banks. Under § 54 of the German Copyright Law where copies are made for other than private use, a fee is payable to the publisher, the amount of the fee to be determined by agreement with the publisher.

© by Springer-Verlag Berlin Heidelberg 1980
Printed in Germany

Printing and binding: Beltz Offsetdruck, Hemsbach/Bergstr.
2145/3140-543210

PREFACE

“Anyone who can spell in English can’t be very bright.”

– George Bernard Shaw

The automatic detection and correction of spelling errors by computers has been a subject of interest for a long time. (Our literature search revealed work as early as 1957.) There have been several papers investigating various algorithms and showing their application to various tasks, generally data entry. Now, however, with the increased interest in computer based text processing (word processing) and the storage of large amounts of textual information in computers (data bases), we suggest that spelling correction will become commonplace. This volume brings together the diverse and scattered work on this topic and shows how it can be applied to create a real general purpose spelling corrector.

The theory behind spelling error detection and correction is only one of the considerations of this volume, however. More generally, we consider the design of the program to implement that theory. This work is an *experiment in program design*. Our objective is to demonstrate the application of modern program design principles on a large real (non-toy) program. The approach taken is:

1. A complete literature search provides the background information upon which our program will be based, introducing the problems and solutions which have been considered before.
2. Using our new familiarity with the subject, a complete, top-down program design will be created.
3. Finally, the design will be implemented by converting it into working code. The complete working program is given in Appendix III.

The purpose of this project is as much to investigate and demonstrate modern programming principles as to produce a spelling corrector. We feel that this project can be used as an example of the design of a large program. As such, it is suitable for study.

Consider using this volume in the classroom. The topic appeals to students; it is new, modern and has the appearance of intelligence. The simplest spelling error detector is an exercise in data structures, but it can be expanded and developed to illustrate many topics (disk I/O, interactive programming, data structures, performance, design, ...). Our experience has shown that it can be a very satisfactory programming project. Combining such a project with the reading of this volume provides both the experience and the example needed to motivate modern programming principles.

The program of Appendix III is available in machine readable form from the author (for a minimal cost to cover the computer time for preparing a copy of the program and shipping charges) for educational purposes.

James L. Peterson
Department of Computer Sciences
The University of Texas at Austin
Austin, Texas 78712

Acknowledgements - Discussions with William B. Ackerman, Robert Amsler, Michael Conner and Leigh Power helped to understand the problems of spelling programs. Jeanne Peterson assisted with proofreading and defining the style of this volume. Art Rinn, Emmett Griner, Mark McCulloch, Barney McCartney and the staff of Texas Student Publications aided in the preparation of the typeset copy. I greatly appreciate their help and time.

TABLE OF CONTENTS

Part I Background

1. Introduction	1
2. Token Lists	3
3. TYPO	5
4. Dictionary Look-up	6
5. Interactive Spelling Checkers	8
6. Spelling Correction	11
7. Affix Analysis	15

Part II Design of a Spelling Program

8. Introduction	18
9. The Main Program	21
10. <i>read_and_obey_speller_directive</i>	24
11. Command Input	26
12. <i>read_and_obey_speller_directive</i> (continued)	31
13. <i>check_spelling</i>	32

14. <i>get_token</i>	34
15. Dictionary Structure and Search	51
16. What To Do With Misspelled Tokens	72
17. What To Do With Misspelled Tokens (continued)	79
18. Local Dictionaries	87
19. User Control of the Speller	104
20. Memory Management	112
21. Implementation	117
22. Improvements	119
23. Conclusions	122
 Bibliography	 123
 Appendix I – The 258 Most Common English Words	 130
 Appendix II – Definition of Common Word Graph	 132
 Appendix III – A Complete Spelling Program	 136
 Index	 207