# The Unicode Cookbook for Linguists

## Managing writing systems using orthography profiles

Steven Moran

Michael Cysouw

language science press

Translation and Multilingual Natural Language Processing

Editors: Oliver Czulo (Universität Leipzig), Silvia Hansen-Schirra (Johannes Gutenberg-Universität Mainz), Stella Neumann (RWTH Aachen), Reinhard Rapp (Johannes Gutenberg-Universität Mainz)

In this series:

# The Unicode Cookbook
# for Linguists

## Managing writing systems using orthography profiles

Steven Moran

Michael Cysouw

Steven Moran , Michael Cysouw. 2018. *The Unicode Cookbook for Linguists*: *Managing writing systems using orthography profiles* (Translation and Multilingual Natural Language Processing 10). Berlin: Language Science Press.

# Preface

This text is meant as a practical guide for linguists and programmers, who work with data in multilingual computational environments.[1] We introduce the basic concepts needed to understand how writing systems and character encodings function, and how they work together.

The intersection of the Unicode Standard and the International Phonetic Alphabet is often met with frustration by users. Nevertheless, the two standards have provided language researchers with a consistent computational architecture needed to process, publish and analyze data from many different languages. We bring to light common, but not always transparent, pitfalls that researchers face when working with Unicode and IPA.

In our daily working lives, we use quantitative methods to compare languages and uncover and clarify their phylogenetic relations. However, the majority of lexical data available from the world's languages is in author- or document-specific orthographies. Having identified and overcome the pitfalls involved in making writing systems and character encodings syntactically and semantically interoperable (to the extent that they can be), we created a suite of open-source Python and R tools to work with languages using profiles that adequately describe their orthographic conventions. Using orthography profiles and these tools allows users to segment text, analyze it, identify errors, and to transform it into different written forms.

We welcome comments and corrections of this text, our source code, and the use case studies that are supplement to this book.[2] Please use the issue tracker, email us directly, or find the book on PaperHive, where readers can leave questions and comments.[3]

Steven Moran <bambooforest@gmail.com>
Michael Cysouw <cysouw@mac.com>

---

[1]https://github.com/unicode-cookbook/cookbook
[2]https://github.com/unicode-cookbook/recipes
[3]https://paperhive.org/

# Acknowledgments

# Contents

# Did you like this book?

This book was brought to you for free

# Change backtitle in localmetadata.tex

Change backbody in localmetadata.tex