# THE LINUX PROGRAMMING INTERFACE

A Linux and UNIX® System Programming Handbook

## MICHAEL KERRISK

# THE LINUX PROGRAMMING INTERFACE

A Linux and UNIX® System Programming Handbook

**MICHAEL KERRISK**

This logo applies only to the text stock.

# BRIEF CONTENTS

# CONTENTS IN DETAIL