M Beck, H Böhme, M Dziadzka,
U Kunitz, R Magnus, C Schröter, D Verworner

# LINUX
# KERNEL
# PROGRAMMING

## THIRD EDITION

CD-ROM INCLUDED
Now revised and updated
to cover Linux 2.4

# CONTENTS

CHAPTER 10    MULTIPROCESSING 276

APPENDIX A    SYSTEM CALLS 284

APPENDIX B    KERNEL-RELATED COMMANDS 373