Heike Hänlein

# Studies in authorship recognition – a corpus-based approach

# Contents