# Abstract

This report discuss the methodology and implementation of two significant module segmentation & recognition part of Optical Character recognition system for scanned image. This methodology consists of three steps: At first step a segmentation approach is used in order to detect text lines, words and characters. In this step the most decisive part was to segment the characters whose regions are overlapped. We have solved this problem with a new approach we called curve segmentation. In second step these segmented characters are then used for recognition to determine the expected sequence of characters. We extracted some features for Bangla characters then examine a simple pattern-recognition system designed and simulated using KNN.

# Acknowledgments

Throughout this thesis, we have received support from many people whom we wish to acknowledge here. First and foremost, we are deeply indebted to our principal supervisor Dr. Muhammed Zafar Iqbal for his advice and for being a ceaseless motivator during the last 6 month. He has always encouraged us to keep an open mind and to push the envelope at all times. We would also like to show our gratitude to our co-supervisor Md. Sabir Ismail & Md. Saiful Islam for sharing their pearls of wisdom with us during the course of this research. And we are so grateful to our supervisor & co-supervisors for arranging special meeting in every week and give their valuable time to us, which have boosted our workflow.

# TABLE OF CONTENTS

**Page**

# List of Tables

# List of Figures