# ACCELERATED SMPP DECODER IMPLEMENTATION BASED ON GPU

Prabath Sanjeewa Abeygunawardana Weerasinghe

Registration Number : 148244C

Degree of Master of Science

Department of Computer Science and Engineering

University of Moratuwa

Sri Lanka

December 2017

# Declaration

I declare that this is my own work and this dissertation does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works.

Candidate

....................................                                                    ..............................

**Prabath Sanjeewa Abeygunawardana Weerasinghe**                    **Date**

I certify that the declaration above by the candidate is true to the best of my knowledge and he has carried out research for the Masters Dissertation under my supervision.

....................................                                                    ..............................

**Prof. Sanath Jayasena**                                                          **Date**

# Acknowledgment

This project would not have been possible without the support of many people.

A special thank goes to my supervisor, Prof. Sanath Jayasena for his valuable guidance and immense support. Many thanks to our MSc research project coordinators, Dr. Shehan Perera and Dr. Malaka Walpola, for their dedication and support. Thanks to all the lecturers at the Faculty of Computer Science and Engineering, University of Moratuwa, for their valuable advices.

# Abstract

Graphic processing unit (GPU) provides a low-cost but powerful hardware platform for implementing massively parallel high performance systems. The capabilities of GPUs have been used to provide fast and low cost solutions in areas such as machine learning, complex simulations such as global warming and genetic engineering and network traffic processing.

Our research is focused on using GPUs to accelerate the decoding process of the Short Message Peer to Peer (SMPP) protocol. SMPP protocol is used to exchange *Short Messages* between Short Message Service Centers (SMSCs) and TCP/IP based applications. From the point of view of a user, a SMS taking few seconds is acceptable and therefore a SMSC is mostly focused on achieving a higher throughput than a low per packet latency.

We have developed a SMPP decoder library in C with GPU support. It supports both CPU based and GPU based decoding. The library also includes two primary APIs. The first API is for general usage by any C based application and the second API could be used by Java Native Interface (JNI) based application.

We have evaluated the performance of the library in both CPU and GPU modes and compared it with a SMPP server based on Cloudhopper, a Java implementation of SMPP protocol. The evaluation shows around five times throughput gain in GPU mode over both Java and C based CPU modes.

# Contents

# List of Figures

# List of Abbreviations

GPU         Graphics Processing Unit

SMPP        Short Message Peer to Peer

CPU         Central Processing Unit

SMS         Short Message System

SMSC        Short Message Service Center

ESME        External Short Message Entity

PDU         Protocol Data Unit

TLV         Tag Length Value

CUDA        Compute Unified Device Architecture

SM          Streaming Multiprocessor

SP          Streaming Processor

UMA         Unified Memory Access