Encrypted Search: Enabling Standard Information Retrieval Techniques for Several New Secure Index Types While Preserving Confidentiality Against an Adversary With Access to Query Histories and Secure Index Contents

By Alexander R. Towell, Bachelor of Science

A Thesis Submitted in Partial
Fulfillment of the Requirements
for the Degree of
Master of Science
in the field of Computer Science

Advisory Committee:

Hiroshi Fujinoki, Chair

Gunes Ercal

Tim Jacks

Graduate School
Southern Illinois University Edwardsville
September, 2015

ABSTRACT

ENCRYPTED SEARCH: ENABLING STANDARD INFORMATION RETRIEVAL
TECHNIQUES FOR SEVERAL NEW SECURE INDEX TYPES WHILE PRESERVING
CONFIDENTIALITY AGAINST AN ADVERSARY WITH ACCESS TO QUERY
HISTORIES AND SECURE INDEX CONTENTS

by

ALEXANDER R. TOWELL

Chairperson: Professor  Hiroshi Fujinoki

*Encrypted Search* is a way for a client to store searchable documents on untrusted systems
such that the untrusted system can *obliviously* search the documents on the client's behalf,
i.e., the untrusted system does not know what the client is searching for nor what the
documents contain. Several new secure index types are designed, analyzed, and
implemented. We analyze them with respect to several performance measures:
confidentiality, time complexity, space complexity, and search retrieval accuracy. In order to
support rank-ordered search, the secure indexes store frequency and proximity information.
We investigate the risk this additional information poses to confidentiality and explore ways
to mitigate said risk. Separately, we also simulate an adversary who has access to a history
of encrypted queries and design techniques that mitigate the risk posed by this adversary.

KEYWORDS: (Encrypted Search, Information Leaks, Perfect Hash Filter, Query Obfuscation,
Secure Indexes)

# ACKNOWLEDGEMENTS

I would like to recognize my better half, Kimberly Wirts, for her patience and encouragement.

I would like to thank my parents, Bill and Sharon Towell, for their constant support throughout the entire process.

I would like to express my gratitude to my thesis advisor, Prof. Hiroshi Fujinoki, for introducing me to my thesis topic and lending me his time and support.

I would like to thank the rest of my thesis committee, Prof. Gunes Ercal and Prof. Tim Jacks, for their insights and challenging questions.

TABLE OF CONTENTS

LIST OF FIGURES

LIST OF TABLES