Ain Shams University
Faculty of Computer and Information Sciences
Information Systems Department

# Efficient Processing of Continuous Queries based on Cloud Computing

A thesis submitted in partial fulfillment of the requirements for
the degree of MSc in Computer and Information Sciences
To
Information Systems Department,
Faculty of Computer and Information Sciences,
Ain Shams University

**By**

## Fatma Mohamed Mahmoud Najib

Teaching Assistant at Information Systems Departments
Faculty of Computer and Information Sciences, Ain Shams University

**Under Supervision of**

## Prof. Dr. Mohamed Fahmy Tolba

Scientific Computing Department
Faculty of Computer and Information Sciences, Ain Shams University

## Prof. Dr.  Nagwa Lotfy Badr

Information Systems Department
Faculty of Computer and Information Sciences, Ain Shams University

## Assoc. Prof. Rasha Mohamed Ismail

Information Systems Department
Faculty of Computer and Information Sciences, Ain Shams University

September – 2016

# Acknowledgement

First of all thanks to Allah for giving me the opportunity and the capability to finish this thesis. Great thanks to everyone who assists and encourages me to complete this work.

I would like to deeply thank Prof. Dr. Mohamed Fahmy Tolba for his leadership, support and encouragement all time.

Also, I would like to thank Prof. Dr. Nagwa Badr for her supervision, instructions and help throughout this thesis.

I would also like to thank Assoc. Prof. Rasha Ismail for her guidance, assistance and support through the work.

Finally, all thanks to my sweetie mother for being the reason to any success in my life. A special gratitude and love to my father for his great and unfailing support and love. I would also like to thank every member of my family for their encouragement.

# Abstract

*Many recent applications in several domains such as sensor networks and financial applications generate continuous, rapid, and time varying datasets which are called data streams. Data streams require real time processing. In most database systems, the query optimizers select a single plan to process all streams tuples which is not efficient with the streams changeable nature. Also there is a little research effort has been made towards the multiple data streams queries' simultaneous execution. In addition applying data streams' multi-directional optimization over an optimized and elastic environment has not been much considered. Thus in this thesis we proposed combined frameworks and different optimization algorithms to solve these problems.*

*First, we proposed the optimized query mesh for data stream (OQMDS) framework. In which data streams are processed over multiple query plans. Each plan is used to process a cluster of data that have nearest properties. We proposed the Optimized Iterative Improvement Query Mesh (OII-QM) and Non-Search based Query Mesh (NS-QM) algorithms, to efficiently generate the multiple plans. The proposed algorithms improves the optimization time by 70.3%, the execution time by 21.8%, the execution overheads by 80% and the memory usage by 96% over the II-QM algorithm.*

*Then in this thesis the Continuous Query Optimization based on Multiple Plans framework for data streams over the cloud environment (CQOMP) was proposed. CQOMP provides an optimized streams processing over the cloud. The Optimized Multiple plans (OMP) and the Auto Scaling Cloud Query Mesh (AS-CQM) algorithms were proposed for streams processing over multiple query plans on cloud computing. The proposed OMP improves the performance in terms of the execution time by 83.5%, 47.7%, and the throughput by 69.7%, 40% over the operator-set-cloud methodology and the NS-QM algorithm. The elastic configurations of the proposed AS-CQM increases utilizing cloud processing resources by 33.8% and reduces the costs by 50% over the static configuration.*

*Finally in this thesis the multiple queries optimization based on partitioning (MQOP) framework was proposed to efficiently execute multiple queries simultaneously on the cloud environment. The optimized global plan (OGP) and the optimized global plan based on partitioning (OGPP) algorithms were proposed for jointly executing multiple continuous queries over an optimized global plan to each cluster of data on the cloud. The proposed OGP improves the execution time by 80% and the throughput by 76.8% over the operator tree technique. The proposed OGPP algorithm improves the performance in terms of execution time by 61.1%, 72.6%, and the throughput by 55.5%, 66.5% over the compile time optimization method and the operator tree technique.*

# Table of Contents

**on Cloud Environment**

# List of Figures

# List of Tables

# List of Algorithms

# List of Abbreviations

| | |
|---|---|
| AMR | The adaptive multi-route query processing system |
| AS-CQM | The proposed Auto Scaling Cloud Query Mesh algorithm |
| A-SEGO | The Adaptive Sharing based Extended Greedy Optimization Approach |
| CANN query | The continuous aggregate nearest neighbor query |
| CEP | The Complex Event Processing framework |
| CQOMP | The proposed Continuous Query Optimization based on Multiple Plans framework for data streams over the cloud environment |
| CRNN query | The continuous range nearest neighbor query |
| DSPS | Distributed Stream Processing Systems |
| EM | The Expectation Maximization algorithm |
| IEVA | The Improved Event Based Algorithm |
| II-QM | The Iterative Improvement Query Mesh algorithm |
| kNN query | The k nearest neighbor query |
| M3 | The Main-Memory MapReduce |
| MkNN query | The mobile k nearest neighbor query |
| MQOP | the proposed multiple queries optimization based on partitioning framework |
| NS-QM | The proposed Non-Search based Query Mesh algorithm |
| OGP | The proposed optimized global plan algorithm |
| OGPP | The proposed optimized global plan based on partitioning algorithm |
| OII-QM | The proposed Optimized Iterative Improvement Query Mesh algorithm |
| OMP | The proposed Optimized Multiple plans algorithm |
| OQMDS | The proposed optimized query mesh for data stream framework |
| PDAMS | The power and deadline-aware multicore scheduling algorithm |
| QM | The Query Mesh framework |
| RLD | Robust Load Distribution algorithm |
| SQO | The semantic query optimization approach |
| VGQ-Vor | The virtual grid quadtree with Voronoi diagram |
| Vk NN query | The visible k nearest neighbor query |
| VM | Virtual machine |

# Chapter 1

---

# Introduction

---

# Chapter 1    Introduction

## 1.1 Overview

Most of the recent applications such as sensor networks [1-4], traffic analysis applications [5], and tracking moving objects [6] generate a continuous, rapid, time varying and dynamic data element that are called data streams [7-12]. These applications present new challenges which are not handled by traditional techniques. Thus continuous queries are considered to process data streams which have a continuous and dynamic nature. Continuous queries are evaluated continuously with the continuous arrival of data streams over time [13-16].

Traditional database systems uses optimizers that process all data stream tuples based on a single query plan. This plan is not efficient with data streams that have continuous variations over time. For example, a continuous query with three operators (op1, op2, op3) and the query result is updated every five minutes for the latest data streams tuples over this interval. Because of the continuous nature of data streams over time, it is not efficient to process all tuples over all runs (each five minute) or even in the same run using one query plan. Consider that the best query plan to process the tuples within the first two minutes is (op1 - op2 – op3) but the best plan for the tuples in the next minute is (op2 - op3- op1) and so on. Therefore data streams processing using multiple query plans is a good alternative solution. Each query plan is the most suitable plan to process a cluster of tuples that have the nearest characteristics.