# Advances in Cluster Editing: Linear FPT Kernels and Comparative Implementations

**Peter E. Shaw**
**B. Comp. Sci. (Hons I )**

*A thesis submitted in fulfilment
of the requirements for the degree of*

**Doctor of Philosophy**

School of Electrical Engineering
and Computer Science

The University of Newcastle
Callaghan N.S.W. 2308
Australia

February, 2010

*The* UNIVERSITY
*of* NEWCASTLE
AUSTRALIA

**STATEMENT OF ORIGINALITY**

The thesis contains no material which has been accepted for the award of any other degree of diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. I give consent to this copy of my thesis, when deposited in the University Library, being made available for loan and photocopying subject to the provisions of the Copyright Act 1968.

I hereby certify that the work embodied in this thesis is the result of original research and has not been submitted for a higher degree to any other University or Institution.

I hereby certify that parts of the work embodied in this thesis have been done in collaboration with other researchers. I have included as part of the thesis a statement clearly outlining the extent of collaboration, with whom and under what auspices.

I hereby certify that the work embodied in this thesis contains a number of published scholarly work (conference proceedings) of which I am a joint author. I have included a written statement, attesting to my contribution towards the joint publications/scholarly work.

<div align="right">Peter E. Shaw</div>

# Acknowledgements

*To my darling wife Eunmee Chang-Shaw whom I love*
*so much.*

Rejoice in the Lord always. I will say it again: Rejoice! Let your gentleness be evident to all. The Lord is near. Do not be anxious about anything, but in everything, by prayer and petition, with thanksgiving, present your requests to God. And the peace of God, which transcends all understanding, will guard your hearts and your minds in Christ Jesus.

*Philippians 4:4–7 (New International Version)*

# Contents

# Abstract

Experience has shown that clustering objects into groups is a useful way to analyze and order information. It turns out that many clustering problems are intractable. Several heuristic and approximation algorithms exist, however in many applications what is desired is an optimum solution. Finding an optimum result for the Cluster Editing problem has proven non-trivial, as Cluster Editing is $\mathcal{NP}$-Hard [KM86], and $\mathcal{APX}$-Hard, and therefore cannot be approximated within a factor of $(1 + \epsilon)$ unless $\mathcal{Poly} = \mathcal{NP}$ [SST04].

The algorithmic technique of Parameterized Complexity has proven an effective tool to address hard problems. Recent publications have shown that the Cluster Editing problem is Fixed Parameter Tractable ($\mathcal{FPT}$). That is, there is a fixed parameter algorithm that can be used to solve the Cluster Editing problem.

Traditionally, algorithms, in computer science, are evaluated in terms of the time needed to determine the output as a function of input size only. However, typically in science most real datum contains inherent structure. For Fixed Parameter Tractable (FPT) algorithms, permitting one or more parameters to be given in the input, to further define the question, allows the algorithm to take advantage of any inherit structure in the data [ECFLR05].

A key concept of FPT is kernelization; that is, reducing a problem instance to a core hard sub-problem. The previous best kernelization technique for Cluster Editing was able to reduce the input to within $k^2$ [GGHN05] vertices, when parameterized by $k$, the edit distance. The edit distance is the number of edit operations required to transform the input graph into a cluster graph (a disjoint union of cliques). Experimental comparisons in [DLL$^+$06] showed that significant improvements were obtained using this reduction rule for the Cluster Editing problem. The study reported in this thesis presents three polynomial-time, many-to-one kernelization algorithms for the Cluster Editing problem, the best of these algorithms produces a linear kernel of at most $6k$ vertices.

In this thesis, we discuss how using new $\mathcal{FPT}$ techniques including *extremal method compression routine* and *modelled crown* reductions [DFRS04] can be used to kernelize the input for the Cluster Editing problem. Using these new kernelization techniques, it has been possible to improve the number of vertices in the data sets that can be solved optimally, from the previous maximum of around 150 vertices to over 900. More importantly, the edit distance of the graphs that could be solved as also increased from around $k = 40$ to more than $k = 400$.

This study also provides a comparison of three inductive algorithmic techniques:

- *compression routine using a constant factor approximation* – Compression Crown Rule Search Algorithm;

- *extremal method* (coordinatized kernel) [PR05], using a *constructive* form of the boundary lemma – Greedy Crown Rule Search Algorithm;

- *extremal method*, using an auxiliary (TWIN) graph structure – Crown Rule TWIN Search Algorithm.

Algorithms derived using each of the above techniques to obtain linear kernels for the Cluster Editing problem have been evaluated using a variety of data with different exploratory properties. Comparisons have been made in terms of reduction in kernel size, lower bounds obtained and execution time.

Novel solutions have been required to obtain approximations within a reasonable time, for the Cluster Editing problem that is within a factor of four of the edit distance (minimum solution size). Most approximation methods performed very badly for some graphs and well for others. Without any guide regarding the quality of the result, a very bad result can be assumed to be close to optimum.

Our study has found that just using the highest available lower bound for the approximation is insufficient to improve the result. However, by combining both the highest lower bound obtained and the reduction obtained using kernelization, a 30-fold improvement in the approximation performance ratio is achieved.

# List of Figures

# List of Tables