

Algorithm Selection for Maximum Common Subgraph

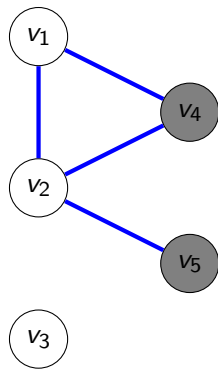
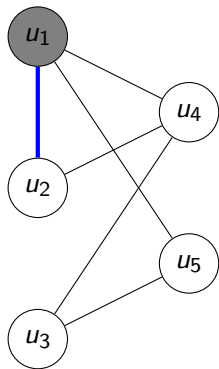
Paulius Dilkas

Supervisors: Dr Patrick Prosser and Dr Ciaran McCreesh

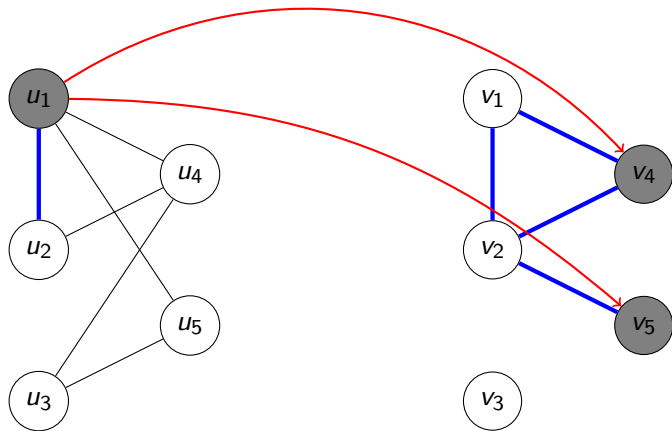
School of Computing Science
University of Glasgow

23rd March 2018

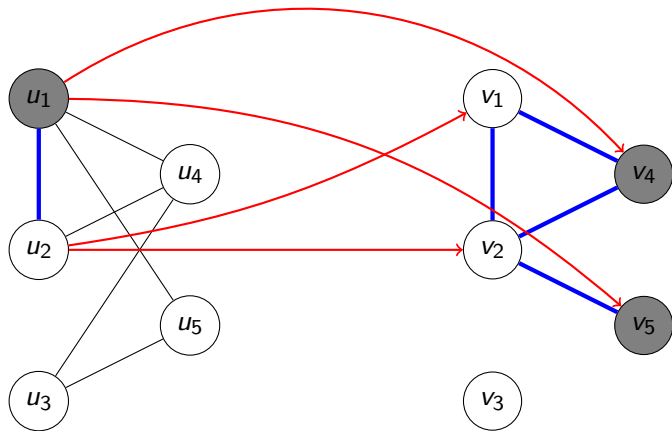
The Problem: Maximum Common Subgraph



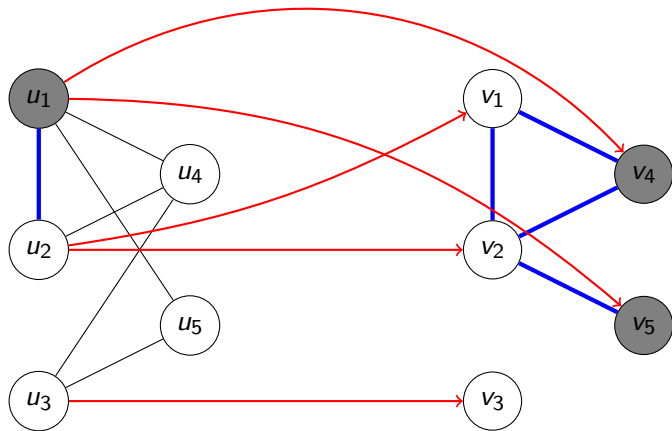
The Problem: Maximum Common Subgraph



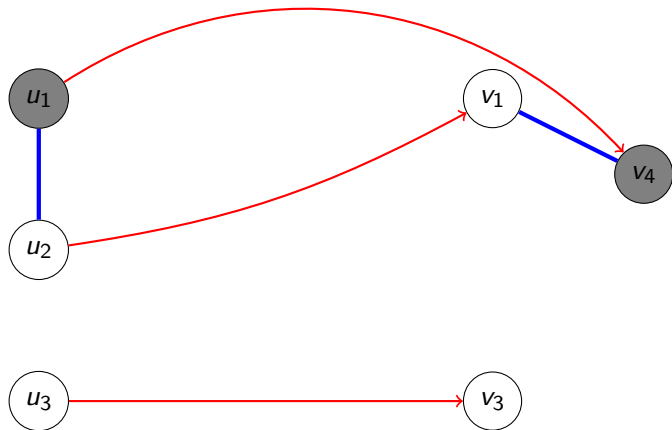
The Problem: Maximum Common Subgraph



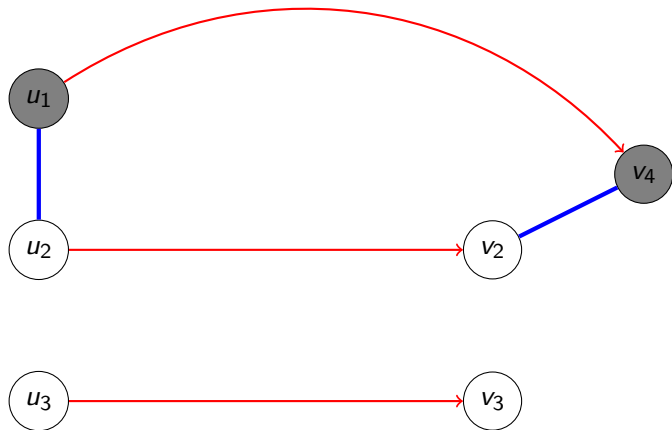
The Problem: Maximum Common Subgraph



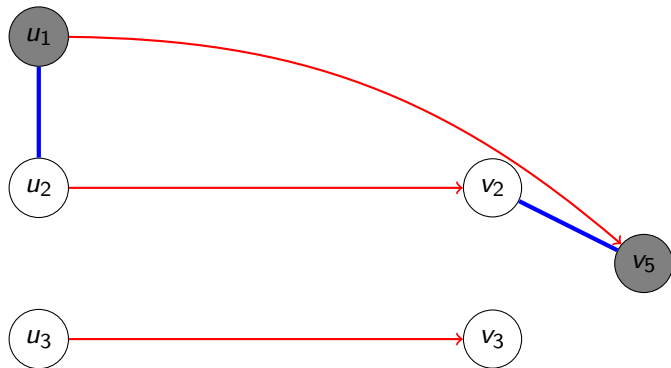
The Problem: Maximum Common Subgraph



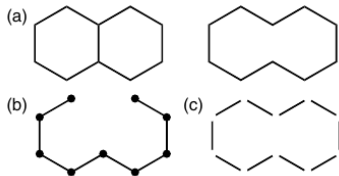
The Problem: Maximum Common Subgraph



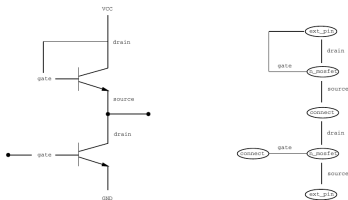
The Problem: Maximum Common Subgraph



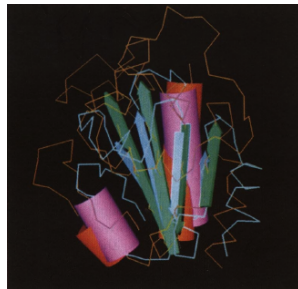
Why Is It Important?



Source: Ehrlich and Rarey 2011



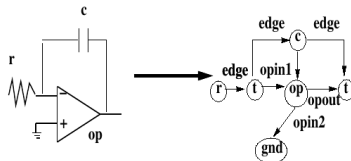
Source: Cook and Holder 1994



Source: M. Grindley et al. 1993

circuit

graph representation



Source: Djoko, Cook and Holder 1997

Existing Ways to Solve It

- ▶ MCSPLIT, MCSPLIT↓
 - ▶ McCreesh, Prosser and Trimble 2017
- ▶ clique encoding
 - ▶ McCreesh, Ndiaye et al. 2016
- ▶ k ↓
 - ▶ Hoffmann, McCreesh and Reilly 2017

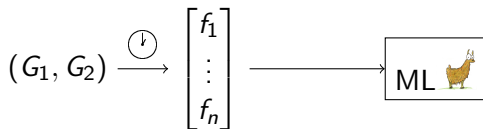
Solution: Algorithm Selection

(G_1, G_2)

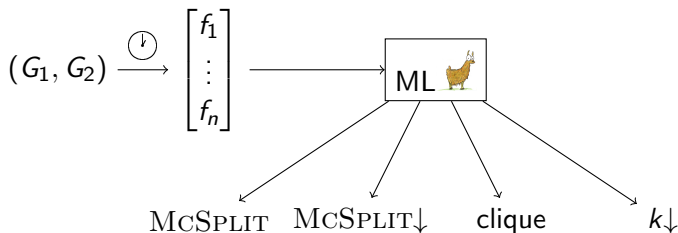
Solution: Algorithm Selection

$$(G_1, G_2) \xrightarrow{\text{⌚}} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

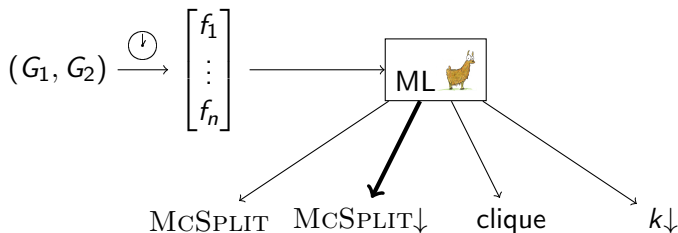
Solution: Algorithm Selection



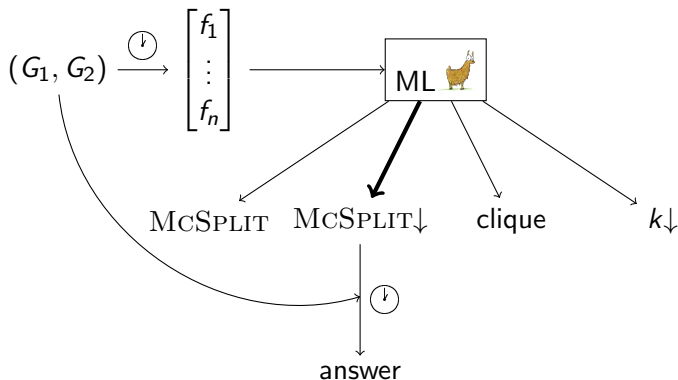
Solution: Algorithm Selection



Solution: Algorithm Selection



Solution: Algorithm Selection



The Process

- ▶ Three cases:
 - ▶ no labels
 - ▶ vertex labels
 - ▶ vertex and edge labels

The Process

- ▶ Three cases:
 - ▶ no labels
 - ▶ vertex labels
 - ▶ vertex and edge labels
- ▶ How many different labels?
 - ▶ Measured as a percentage of the number of vertices/edges
 - ▶ 5%, 10%, 15%, 20%, 25%, 33%, 50%

The Process

- ▶ Three cases:
 - ▶ no labels
 - ▶ vertex labels
 - ▶ vertex and edge labels
- ▶ How many different labels?
 - ▶ Measured as a percentage of the number of vertices/edges
 - ▶ 5%, 10%, 15%, 20%, 25%, 33%, 50%
- ▶ Run every algorithm on every instance
 - ▶ $\sim 500,000$ experiments for each algorithm

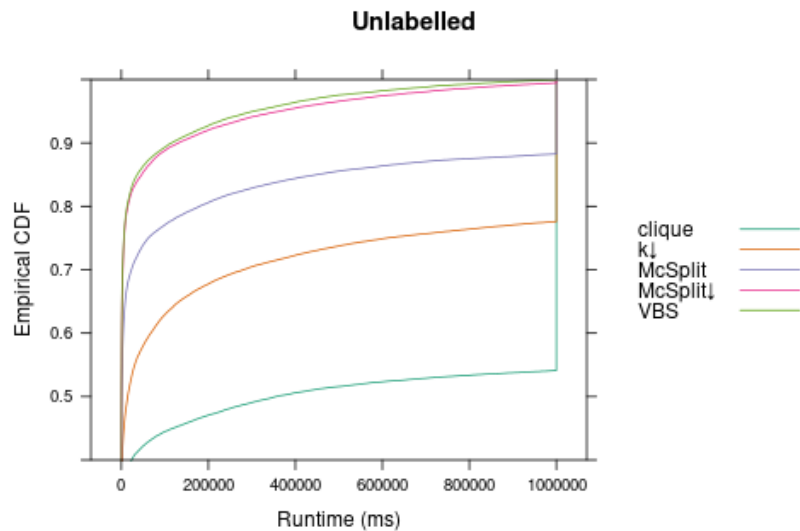
The Process

- ▶ Three cases:
 - ▶ no labels
 - ▶ vertex labels
 - ▶ vertex and edge labels
- ▶ How many different labels?
 - ▶ Measured as a percentage of the number of vertices/edges
 - ▶ 5%, 10%, 15%, 20%, 25%, 33%, 50%
- ▶ Run every algorithm on every instance
 - ▶ $\sim 500,000$ experiments for each algorithm
- ▶ Identify and extract features
 - ▶ 34 in total
 - ▶ some from Kotthoff, McCreesh and Solnon 2016
 - ▶ some new

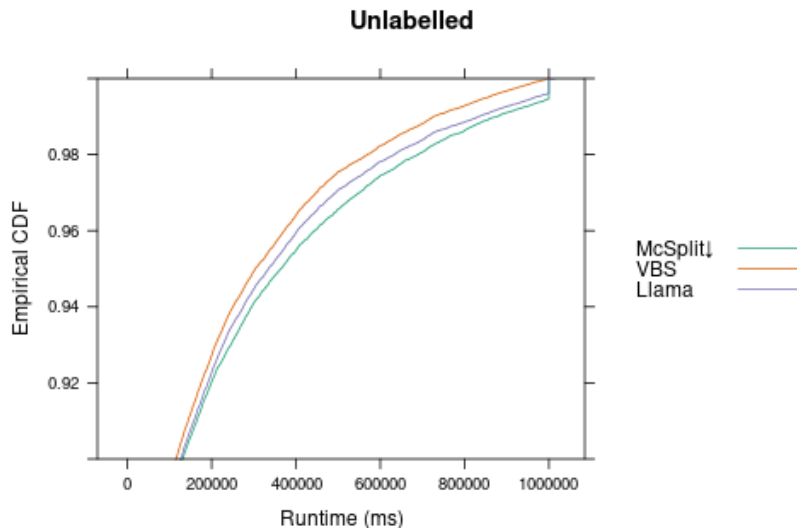
The Process

- ▶ Three cases:
 - ▶ no labels
 - ▶ vertex labels
 - ▶ vertex and edge labels
- ▶ How many different labels?
 - ▶ Measured as a percentage of the number of vertices/edges
 - ▶ 5%, 10%, 15%, 20%, 25%, 33%, 50%
- ▶ Run every algorithm on every instance
 - ▶ $\sim 500,000$ experiments for each algorithm
- ▶ Identify and extract features
 - ▶ 34 in total
 - ▶ some from Kotthoff, McCreesh and Solnon 2016
 - ▶ some new
- ▶ Train machine learning models
- ▶ Evaluate their performance and usefulness

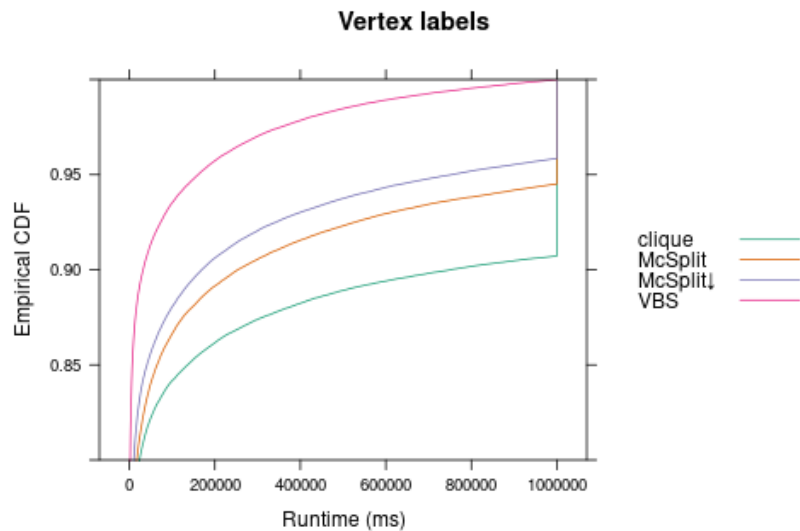
Results



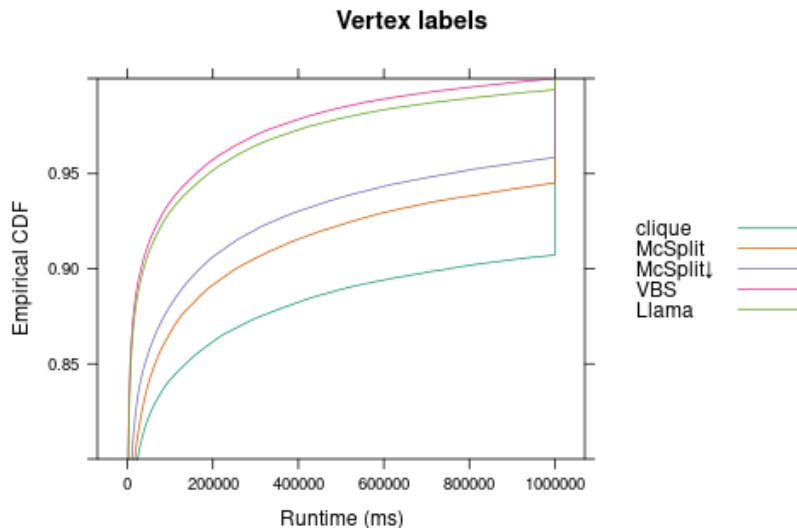
Results (27%)



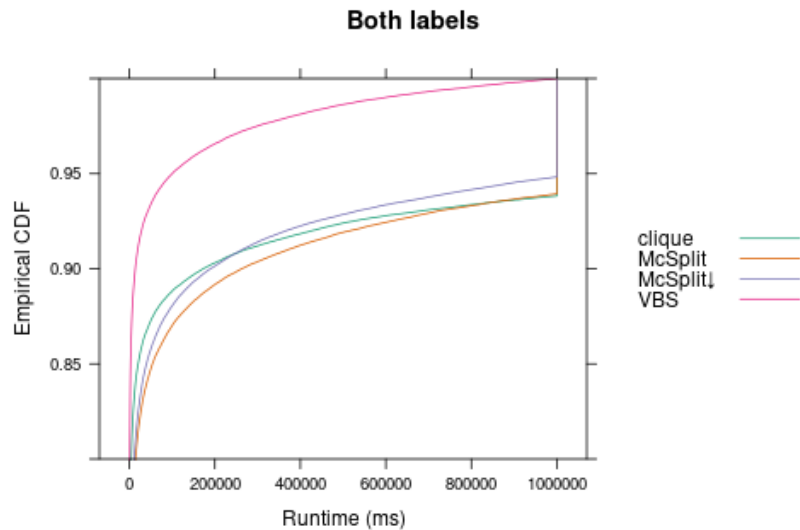
Results



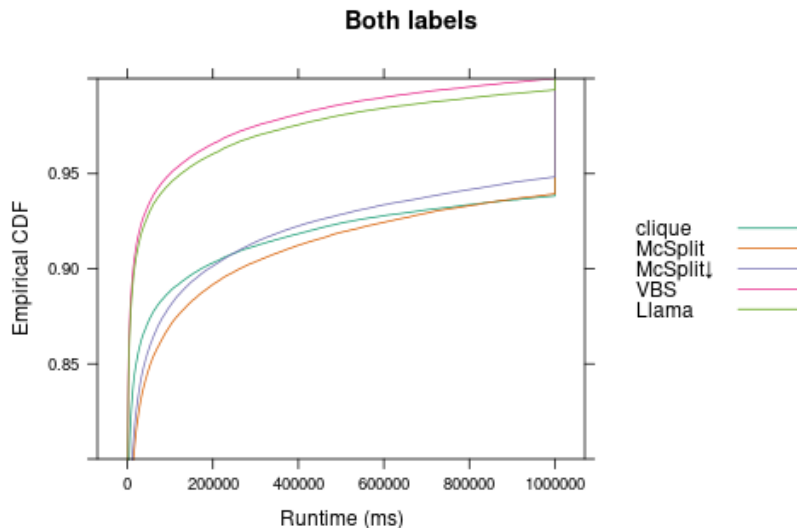
Results (86%)



Results



Results (88%)



Other Accomplishments

- ▶ Identified important features:
 - ▶ labelling percentage, standard deviation of degrees
- ▶ Discovered how algorithms' performance changes with fewer labels
- ▶ Extended k_{\downarrow} to support vertex labels
 - ▶ using neighbourhood degree sequences
- ▶ Defined and developed new algorithms capable of switching between MCSPLIT and the clique encoding