

Final Report — Optimization of Subgraph Isomorphism Algorithm

by Nelson Batista, Max Inciong, Francesca Truncale

Senior Project II

Professor Jianting Zhang

Fall 2017

Contents

1	Introduction	1
1.1	Background	1
1.2	Motivation	1
1.3	Data Set	2
2	Ullmann Algorithm	2
2.1	Naive Approach	2
2.2	“Refine M” Procedure	3
3	Alternative Algorithms	3
3.1	VF2 Algorithm	3
3.2	RI Algorithm	6
4	Implementation of the Algorithm	7
4.1	Python Implementation	7
4.2	Initial C++ Implementation	8
4.3	Improved C++ Implementation	8
4.4	Options for Parallelization	8
4.4.1	GPGPU	8
4.4.2	MPI	9
4.4.3	SIMD	9
4.5	Optimizations through OpenMP	10
4.6	Sample Run	10
5	Performance	16
5.1	Python Implementation	16
5.2	C++ Implementation	16
5.2.1	Performance Tradeoffs	16
5.2.2	Results	17
5.3	OpenMP Optimization	17
6	Conclusion	18

7	Appendix	19
7.1	Python Code	19
7.1.1	Graph Class	19
7.1.2	Subgraph Isomorphism Algorithm	21
7.1.3	Main	24
7.2	C++ Code	28
7.2.1	Graph Class	28
7.2.2	Subgraph Isomorphism Algorithm	32
7.2.3	Main	39
7.3	OpenMP Code	47
7.3.1	Subgraph Isomorphism Algorithm	48
7.4	Code Used to Check Performance of C++ code	55
	References	63

1 Introduction

The goal of this project is to understand and optimize an algorithm for solving the NP-complete problem of subgraph isomorphism.

1.1 Background

An *isomorphism* between two graphs is a mapping from the vertices of one graph, say G , to another graph, say H , such that any two vertices which are adjacent in G are mapped to vertices which are adjacent in H . The subgraph isomorphism problem, in turn, is the problem of determining whether an isomorphism exists between a graph G and any of the subgraphs of a larger search graph, H . The issue that makes the subgraph isomorphism problem so difficult to solve for a particular pair of graphs is cycling through the main graph and expanding each node within the graph and checking if each of the generated subgraphs formed from each node is isomorphic to the control graph.

Fortunately, many attempts exist to solve the subgraph isomorphism problem. Among the earliest of these is the famous *Ullmann Algorithm*, proposed by Julian R. Ullmann in his 1976 paper, *An Algorithm for Subgraph Isomorphism*. He first describes a basic, naive approach to finding a mapping from the vertices of a graph to a subgraph of a larger graph. He goes on to define a “refine” procedure to dramatically reduce the number of possible mappings that must be checked.[6] These will be covered in greater detail in the next section.

The goal of our project was to take this algorithm and improve its performance. Subgraph isomorphism is an expensive problem to solve, and finding multiple possible isomorphisms from one graph to subgraphs of another is even more expensive, so even slight improvements in the algorithm’s performance is likely to lead to large performance improvements in any large-scale program which needs to use it often.

1.2 Motivation

Our reasoning for choosing this topic for our project is an apparent lack of especially efficient algorithms for solving the subgraph isomorphism problem. Solutions to the subgraph isomorphism problem are often used to detect similarities in chemical compounds,[6] which may shed light on some of their properties. If this process needs to be done, for example, on a very large database of chemical compounds, or less frequently for several very large compounds, the process may take a very long time to complete. It would be very beneficial to improve the performance of this algorithm, even if only slightly, so that these sorts of use cases can still be handled in a more reasonable amount of time.

1.3 Data Set

The primary graph used to test the algorithm and its performance is an undirected graph consisting of “friends lists” from the social media website, Facebook. Each node of the graph represents a unique user, and, if two nodes are adjacent, then their corresponding users are friends on Facebook.[4]

The graph contains 4,039 nodes and 88,234 edges. The relatively large size of the graph makes it well-suited to use for collection of performance data. As we shall see, the Ullmann algorithm takes a rather long time to complete with a graph of this size, allowing any optimizations and/or parallelizations to be readily apparent.

2 Ullmann Algorithm

2.1 Naive Approach

The paper begins, as many do, with a definition of terms that will be used throughout. We have repeated them here for the sake of clarity later on. We begin by noting that the algorithm models the search for a possible mapping by having all possible mappings as nodes on a tree, on which a depth-first search is performed to find a solution.

The subgraph isomorphism problem is defined as the problem of finding all isomorphisms between a graph $G_\alpha = (V_\alpha, E_\alpha)$ and subgraphs of another graph $G_\beta = (V_\beta, E_\beta)$, with (V_α, E_α) and (V_β, E_β) being the set of vertices and edges of G_α and G_β , respectively. The number of vertices and edges of G_α and G_β , respectively, are (p_α, q_α) and (p_β, q_β) .

The *adjacency matrix* for graph G_α , $[a_{ij}]$, is defined by:

$$a_{ij} = \begin{cases} 1 & \text{if } i \neq j \text{ and } i, j \text{ share an edge} \\ 0 & \text{otherwise} \end{cases}$$

The adjacency matrix for graph G_β , $[b_{ij}]$, is defined equivalently.

Notably, we have the $p_\alpha \times p_\beta$ matrix M , which is a binary matrix of assignments (mappings) from V_α to V_β . That is, if m_{ij} is 1, then vertex i in graph G_α could possibly be mapped to vertex j in G_β . Otherwise m_{ij} is 0. M has a few interesting properties as a result, which will soon be apparent.

Finally, we come to d , which is simply algorithm’s current depth in the search tree. The algorithm starts at $d = 0$, and terminates when $d = p_\alpha$. The matrix M at a particular depth d is denoted as M_d , with solutions to the working instance of the subgraph isomorphism problem being the set of assignments from V_α to V_β , contained within matrices M_{p_α} . M_0 is generated by the following formula:

$$m_{ij} = \begin{cases} 1 & \text{if } \text{degree}(j \in G_\beta) \geq \text{degree}(i \in G_\alpha) \\ 0 & \text{otherwise} \end{cases}$$

Each step of the computation consists of setting all but one of the values of one of the rows of M to 0, and then checking whether exactly one 1 exists in each row. If it does, we have found an isomorphism, since each vertex of G_α has been mapped to exactly one vertex of G_β . If it does not, then we continue onto another row of M and check again. If we have reached the point where we cannot perform the operation on a row of M without violating the condition that each *column* of M contains at most a single 1 (a vertex of G_β cannot be mapped from multiple vertices of G_α), then we backtrack to a previous version of M and try to change the column that we set to 1 in the current row. After exhausting all columns, we conclude there cannot be an isomorphism, since at least one vertex of G_α cannot be mapped to any vertex of G_β . [6]

2.2 “Refine M” Procedure

In the naive algorithm, the entire tree is searched, including branches which cannot possibly contain a solution at any depth. This is a large number of M matrices to evaluate, so Ullmann introduces a “refine M” procedure to reduce the number that must be checked by eliminating nodes of the tree, and thus all of their child nodes, which cannot contain a solution.

The procedure is actually quite simple. We iterate through each 1 m_{ij} in M and check whether, for each neighbor of the vertex in G_α corresponding to i , there exists at least one possible assignment from it to a neighbor of the vertex in G_β corresponding to j . If this condition does not hold, the 1 is changed to a 0. Since doing so may render other 1s invalid, we repeat the refine procedure on the modified M until running the procedure does not change it any further, indicating that the refinements do not produce any new invalid 1s. This obviously involves several layers of nested loops, adding considerable overhead to the overall algorithm. However, this overhead is very little considering the number of possibilities that are subsequently eliminated, making it a very worthwhile operation to perform regardless.

This refine M procedure is run both on the initial M as well as each of M_d , greatly reducing how many iterations we must perform to see that a particular mapping is invalid.

3 Alternative Algorithms

3.1 VF2 Algorithm

One algorithm that we considered using was the VF2 Algorithm, which was developed in 2004 by Luigi P. Cordella et al. We have included a general idea of how the algorithm works here, but greater details can be found in the paper

itself, listed under our references.

Given two graphs $G_1(n_1, b_1)$ and $G_2(n_2, b_2)$ and Mapping $M : n_1 \times n_2$, a subset-mapping solution $M(s)$ takes two subgraphs $G_1(s)$ and $G_2(s)$, obtained by selecting from G_1 and G_2 . only the nodes included in the components of $M(s)$, and the branches connecting them. The following algorithm is defined by the paper to implement VF2.

```

PROCEDURE Match(s)
  INPUT:  an intermediate state s; the initial state  $s_0$  has  $M(s_0)=\emptyset$ 
  OUTPUT: the mappings between the two graphs

  IF  $M(s)$  covers all the nodes of  $G_2$  THEN
    OUTPUT  $M(s)$ 
  ELSE
    Compute the set  $P(s)$  of the pairs candidate for inclusion in  $M(s)$ 
    FOREACH  $p$  in  $P(s)$ 
      IF the feasibility rules succeed for the inclusion of  $p$  in  $M(s)$  THEN
        Compute the state  $s'$  obtained by adding  $p$  to  $M(s)$ 
        CALL Match( $s'$ )
      END IF
    END FOREACH
    Restore data structures
  END IF
END PROCEDURE Match

```

Figure 3.1.1: Description of the VF2 Algorithm.

For each mapping (m, n) , VF2 applies five feasibility rules to the current state and the pair to be added. These rules are denoted as: R_{pred} , R_{succ} , R_{in} , R_{out} , and R_{new} . Only if all feasibility rules are satisfied is (n, m) added to the set of mappings $M(s)$. The general form of the feasibility function is defined as:

$$F_{syn}(s, n, m) = R_{pred} \wedge R_{succ} \wedge R_{in} \wedge R_{out} \wedge R_{new}$$

R_{pred} and R_{succ} represent the predecessors and successors respectively of a given node n in a graph. The other three feasibility rules are meant to prune the tree. R_{in} and R_{out} are both 1-step ahead look ahead while R_{new} is a 2-step look ahead. “Look ahead” refers to the k -look-ahead rules defined in the paper. They are used to reduce the number of states generated, by checking in advance if a consistent state s has no consistent successors after k steps.

The feasibility rules are defined as follows:

$$\begin{aligned}
R_{\text{pred}}(s, n, m) &\Leftrightarrow \\
&(\forall n' \in M_1(s) \cap \text{Pred}(G_1, n) \exists m' \in \text{Pred}(G_2, m) \mid (n', m') \in M(s)) \wedge \\
&(\forall m' \in M_2(s) \cap \text{Pred}(G_2, m) \exists n' \in \text{Pred}(G_1, n) \mid (n', m') \in M(s)),
\end{aligned}$$

$$\begin{aligned}
R_{\text{succ}}(s, n, m) &\Leftrightarrow \\
&(\forall n' \in M_1(s) \cap \text{Succ}(G_1, n) \exists m' \in \text{Succ}(G_2, m) \mid (n', m') \in M(s)) \wedge \\
&(\forall m' \in M_2(s) \cap \text{Succ}(G_2, m) \exists n' \in \text{Succ}(G_1, n) \mid (n', m') \in M(s)),
\end{aligned}$$

$$\begin{aligned}
R_{\text{in}}(s, n, m) &\Leftrightarrow \\
&(\text{Card}(\text{Succ}(G_1, n) \cap T_1^{\text{in}}(s)) \geq \text{Card}(\text{Succ}(G_2, m) \cap T_2^{\text{in}}(s))) \wedge \\
&(\text{Card}(\text{Pred}(G_1, n) \cap T_1^{\text{in}}(s)) \geq \text{Card}(\text{Pred}(G_2, m) \cap T_2^{\text{in}}(s))),
\end{aligned}$$

$$\begin{aligned}
R_{\text{out}}(s, n, m) &\Leftrightarrow \\
&(\text{Card}(\text{Succ}(G_1, n) \cap T_1^{\text{out}}(s)) \geq \text{Card}(\text{Succ}(G_2, m) \cap T_2^{\text{out}}(s))) \wedge \\
&(\text{Card}(\text{Pred}(G_1, n) \cap T_1^{\text{out}}(s)) \geq \text{Card}(\text{Pred}(G_2, m) \cap T_2^{\text{out}}(s))),
\end{aligned}$$

$$\begin{aligned}
R_{\text{new}}(s, n, m) &\Leftrightarrow \\
&\text{Card}(\tilde{N}_1(s) \cap \text{Pred}(G_1, n)) \geq \text{Card}(\tilde{N}_2(s) \cap \text{Pred}(G_2, n)) \wedge \\
&\text{Card}(\tilde{N}_1(s) \cap \text{Succ}(G_1, n)) \geq \text{Card}(\tilde{N}_2(s) \cap \text{Succ}(G_2, n)).
\end{aligned}$$

Figure 3.1.2: Description of the feasibility rules.

When working with graphs with more than 200 vertices, VF2 typically performs better than Ullmann's algorithm. The time complexity of VF2 is $O(n^2)$ at best, and $O(n!n)$ at worst. Ullmann's algorithm ranges from $O(n^3)$ to $O(n!n^2)$. [3]

3.2 RI Algorithm

For this algorithm, we need to generate all possible maps between two graphs, then check if any of those maps is a subgraph isomorphism.

The maps can be represented using a *search space tree*, which has a dummy root. Each node represents a possible match between a vertex in graph G (which is referred to as the pattern graph) in the target graph G' .

Unlike other algorithms, RI orders the vertices of the pattern graph in a way that maximizes the chance that a partial path will be removed. The idea is to introduce edge constraints as early as possible, and are modeled only on the pattern graph.

RI is also modeled for directed graphs, which allows for more pruning since (u, v) does not imply (v, u) (where u, v are vertices).

The majority of the work is performed in a greedy algorithm called `GreatestConstraintFirst`, which is used to find a good sequence of vertices based on the number of neighbors a vertex has.

Using the results of `GreatestConstraintFirst`, the following isomorphism rules remove unfeasible paths:

1. Neither u_i nor $M(u_i)$ is already matched in the current path.
2. The matched vertices are compatible, i.e., $lab(u_i) \equiv lab(M(u_i))$.
3. The number of edges connected to $M(u_i)$ in V' is greater than or equal to the number of edges connected to u_i in V . That is $|\{(v', M(u_i)) \in E'\}| \geq |\{(w, u_i) \in E\}|$ and $|\{(M(u_i), v') \in E'\}| \geq |\{(u_i, w) \in E\}|$. In the case of undirected graphs it verifies that $|\langle M(u_i), v' \rangle \in E'| \geq |\langle u_i, w \rangle \in E|$.
4. The constraints deriving from the topology of the pattern graph up to this point in the path are met, $\forall u_i, u_j \in V$ where $0 \leq j \leq i$ $(u_i, u_j) \in E \Rightarrow (M(u_i), M(u_j)) \in E'$ If edges are labeled, then β is defined, we would also verify the compatibility of the edge labels.

Figure 3.2.3: RI Feasibility Rules

The algorithm does not have a dedicated pruning function, which is another factor for the increase in performance when compared against other algorithms (such as VF2); instead, pruning is done through the matching process.[1]

4 Implementation of the Algorithm

The full source code for all versions can be found in the appendix.

4.1 Python Implementation

The Python code is based largely on an implementation of the Ullmann algorithm provided by an answer on the website *StackOverflow*.^[2] It would not run in its original form, so a few adjustments were needed.

```
def find_isomorphism(graph, subgraph):

    assignments = []
    possible_assignments = [[True]*graph.n_vertices() \
                             for i in range(subgraph.n_vertices())]
    if search(graph, subgraph, assignments, possible_assignments):
        return True
    return matches
```

Figure 4.1.1: Function used to determine whether an isomorphism exists.

The main function is `find_isomorphism`, which is depicted above. It calls the function `search`, which performs the actual work of finding an isomorphism.

The `search` function runs the `update_possible_assignments` function. This particular function has been the subject of refining over the course of the project. At the time of the python implementation, it was the most significant bottleneck.

This was mostly due to the usage of the `has_edge` method of the graph class, shown below.

```
def has_edge(self, vert1, vert2):

    """ Checks if edge connecting vert1 and vert2 is in the graph
    """

    #if adjacent, there's an edge
    return ({vert1, vert2} in self.adjacencies)
```

Figure 4.1.2: Code for the `has_edge` method of our graph class.

This method calls the built-in `in` function of Python lists, which has a runtime complexity of $O(n)$.^[5] Thus, the

`has_edge` method of the Python implementation is $O(n)$.

4.2 Initial C++ Implementation

The initial C++ implementation was based entirely on the Python implementation. Indeed, it represented our best attempt to “translate” the Python code to C++. As one might imagine, it was riddled with compilation errors. That said, we did manage to get it to run, and analysis indicated it performed little better than the Python equivalent. We suspected this was because it retained the same issue that the Python implementation had: calling `has_edge`, an $O(n)$ method, countless times. This motivated a significant improvement.

4.3 Improved C++ Implementation

The improved C++ implementation represents a significant milestone in our project, as it not only improved upon the clarity of the isomorphism algorithm, but also resulted in a much cleaner graph class, with edges implemented as an adjacency matrix. This makes the `has_edge` function $O(1)$, as it consists of simply looking up whether the given vertices share an edge by checking the matrix. This makes its repeated use in `refine_possible_assignments` much more acceptable, a fact which is reflected in the performance.

We also have stored vertices in such a way that each vertex is also stored along with its degree, allowing quick and easy lookup of a vertex’s degree.

The improved C++ implementation follows a bit of a loose interpretation of Ullmann’s description of the algorithm. A version which used the `goto` command with labels was considered, to more closely match said description, but was scrapped in favor of a more conventional approach

4.4 Options for Parallelization

4.4.1 GPGPU

When GPUs (graphics processing units) began including programmable shaders and support for floating point arithmetic, they became capable of doing calculations on matrices or vectors. While early GPGPU programming (General Purpose computing on Graphics Processing Units) required formulating problems as a series of graphical operations, there are now several APIs available which allow programmers to program directly on GPUs as they would normally.

One of these APIs is CUDA, which is a proprietary parallel computing platform created by Nvidia. Source code written in CUDA uses C++ syntax, and only Nvidia GPUs are enabled for CUDA. It has the potential to use more than the available cores from a standard CPU, and is only limited by the number of cores that the GPU contains. This ranges from 8 to 5760 (at the time of writing, the Titan Z has 5760 CUDA cores). The main difference between using

CUDA and OpenMP is that CUDA has access to many more threads. OpenMP relies on using the cores of a CPU to do parallel work; as a result, work is sped up by approximately a factor of 2^8 , depending on how many cores a computer has.

A similar API is OpenCL (Open Computing Language), a framework that allows programs to execute across multiple types of processing hardware. Similar to CUDA, source code is written using C-like syntax. Unlike CUDA, it does not require specialized hardware; OpenCL is an open standard, which can be used across AMD, Nvidia, and Intel GPUs, as well as other processors. However, it is notoriously difficult to both set up for use, and to use.

4.4.2 MPI

MPI stands for Message Passing Interface. This utilizes multiple machines in order to parallelize the work. Unfortunately there is a large difference between parallelization using multithreading and MPI. With MPI, the function runs once on each machine, rather than running in multiple threads on a single machine. This means that in order to increase effectiveness, the workload must be balanced among the machines, otherwise some machines might do much more work than other machines, reducing the effective time saved. While load-balancing is a concern for all forms of parallelization, it is especially important when using MPI, since a program running in parallel using MPI is only as fast as the slowest machine, and there is a large amount of overhead in having to send each machine's results across a local network to other machines for further processing. Thus, one cannot simply use MPI to parallelize a program and expect it to perform faster by some predetermined factor. In the case of graph and subgraph isomorphism, we would have to split the work based on which sections of the graph take longer to find mappings for.

Theoretically, it is possible to combine MPI with OpenMP in order to further increase the efficiency. However, our team did not have a network of machines in order to use MPI. If we did, the time it would take to find an isomorphism would decrease by at most a factor based on the number of processors used, multiplied by the number of cores in each processor capable of using OpenMP.

4.4.3 SIMD

SIMD (single instruction, multiple data) uses a series of data types, which can be utilized on CPUs that support them. These data types allow for data level parallelism, similar to OpenMP. Examples of SIMD instructions that utilize these data types are AVX, SSE, and ARM-NEON vector instructions. These instructions allow for operations to be done on these data types on multiple sections at once. For example, when iterating through a loop, we can load four values of an array and add them all at once in a single operation, rather than loading them one at a time into a separate variable to be added. This would allow the program to only iterate one fourth as many times to complete the computation. Other operations can either be between two vectors or between values within the vectors. However,

taking advantage of SIMD programming would have required us to drastically alter the data structures we currently have in place to incorporate SIMD datatypes.

4.5 Optimizations through OpenMP

Ultimately, we decided to settle on using OpenMP. OpenMP utilizes the multiple cores of a computer's processor in order implement multithreading to run related tasks in parallel. We intend to use it where the bottlenecks in our algorithm are to improve the algorithm's performance. The default number of cores for many computers is four. We intend to place the parallelisms where there tend to be a large number iterations, specifically in loops. This means that, theoretically, the program should run four times faster than without multithreading, using our earlier assumption of four processor cores allotted to the algorithm. It is more likely that it will be approximately 3 to 3.5 times faster, since we may not be able to completely parallelize the bottlenecks as intended. Additionally, using OpenMP, as is the case with any parallelization method, comes with a certain amount of overhead in separating the work to be done. There is also the issue of load balancing: making sure that each thread is performing approximately the same amount of computation so that no thread is idling.

The most important place to place the OpenMP is in the `refine_possible_assignments` function. This function is where the majority of the algorithm's computation takes place, as it has several layers of nested loops, each iterating through vectors which may be very large, depending on the sizes of the graphs being examined. Ideally, we want to balance the loads such that each processor core does an equal amount of work, so as to maximize the amount of parallelization done. However, while we may simply parallelize every iterative process, that might not be efficient enough. Ideally, we would further parallelize in nodes of the search tree which have many possible assignments.

4.6 Sample Run

For visualization purposes, we have provided a sample run of the algorithm. For simplicity, only the C++ implementation is seen here. The graphs used are very simple, as this is only for demonstration. The following is the list of edges defining the search graph:

```
1 2
2 3
```

Figure 4.6.1: Search graph used for walkthrough.

The following is the subgraph we will be searching for an instance of in that graph:

Figure 4.6.2: Subgraph being searched for in search graph.

We can place these into text files, named `example.edges` and `sub.edges` for figures 4.6.1 and 4.6.2, respectively. With these filenames, and assuming the program executable is a file in the same directory named `main`, we may run the program from the terminal using the command `./main example.edges sub.edges`. That is, we pass the text file specifying the search graph in as the first argument, and the corresponding file for the subgraph as the second argument.

The program may also be run in interactive mode by passing in the argument `-i` or `--interactive`, in which case the user is prompted to enter the list of edges for each graph manually in the terminal. This allows for manual entering of the graphs in case files are not available for whatever reason, or for testing simple cases that are not worth creating files for. A full list of arguments along with usage information can be viewed by running the program with the `-h` or `--help` arguments.

After the two files are parsed and the corresponding graph objects are created, the program runs `find_isomorphism` with the subgraph and search graph as arguments, which returns a vector of assignments from the subgraph vertices to the search graph vertices, or an empty vector if no isomorphism exists.

Note that, at this point, the search graph consists of a graph with vertices $(1, 2, 3)$, and the subgraph consists of vertices $(10, 20)$.

In `find_isomorphism`, we first set up the matrix of possible assignments M by running `vector < vector<bool> > possible_assignments = create_possible_assignments(sub, graph);`. This creates the matrix by initializing a vector of size p_α , populated with p_α vectors of size p_β . Each of the values in these vectors is initialized to false. This is all done in one line using one of the constructors provided by the C++ vector class. We then use a nested for loop to iterate through each position of the matrix, setting the position to true if the vertex corresponding to the current column in the search graph has a degree greater than or equal to the degree of the vertex corresponding to the current row in the subgraph. Our graph class stores vertices as pairs of integers. The first integer is the unique value associated with the vertex, and the second integer is the degree of the vertex. The degree is incremented every time a pair which contains the vertex is added to the graph. This way, looking up degree of a vertex is trivial. `create_possible_assignments` returns the matrix created.

The next step is to initialize the various variables used to keep track of our place in the computation. The first of these is `assignments_tree`, which is a vector of max size p_α that will contain the matrix of possible assignments at each depth. This is so that we have a way to access the previous matrices in the event that we need to backtrack. The next variable we need is `cols_used`, which keeps track of which columns (corresponding to search graph

vertices) have already been assigned to a row (corresponding to subgraph vertices). `cols_used[i] == true` if column `i` has been assigned. Each search graph vertex can only be assigned to a single subgraph vertex, so it is important to keep track of which columns have already been assigned. The next variable we will need is `col_depth`. `col_depth[i-1] == k` if column `k` was assigned to a row at depth `i`. The need for this will become clear later on in the example, but in simple terms it is used to keep track of what the last column we tried to assign to a row was, so that we do not continually attempt to assign the same column to the same row after determining doing so does not lead to an isomorphism. The final variable is self-explanatory: `depth`. This initialized to 1, and keeps track of our depth in the search tree, if we represent the sets of all possible assignments as a search tree, with isomorphisms at depth p_α if they exist.

Finally, we pass all of these variables along with the subgraph and search graph to the `search` function, which returns true if an isomorphism is found, and false otherwise. At this point, matrix `M` consists of:

$$\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

Since every vertex in the search graph has degree ≥ 1 , every vertex is a possible assignment for every vertex in the subgraph.

The first thing we do in `search` is refine matrix `M`. This is done first, so that, if the refinement fails, we will know right away, and avoid having to do any unnecessary work. We place the call to the function as the conditional to an if statement. If it returns false, we go into the body of the if statement, which contains only the statement `return false;`, indicating that no isomorphism was found.

In the refinement procedure, to which we have passed the matrix of possible assignments and the two graphs, the first noteworthy variable is the boolean `changes_made`, which keeps track of whether or not we have made any changes to the possible assignments matrix. Since eliminating a possible assignment may make other possible assignments invalid, we must run the refinement procedure on the matrix repeatedly to eliminate these newly-invalid assignments. This may occur several times, so we have a variable to keep track of when we can stop. Thus, the entire procedure is wrapped in a while loop, with `changes_made` being the conditional. In the while loop, we first set `changes_made` to false, and we will set it to true if we mark any possible assignment as invalid.

For convenience, here is the current state of possible assignments once again:

$$\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

We first iterate through each of the rows using a simple for loop. That is, we check each subgraph vertex's possible

assignments for validity. The next noteworthy variable is the boolean `no_one`, which is initially true. This variable, which may be confusingly-named, keeps track of whether or not the current row contains no possible assignments. If this is true, we stop the procedure entirely and return false, since it is no longer possible for an isomorphism between the two graphs to exist. Thus, continuing to prune possible assignments is a pointless effort if we have already found there are no possible assignments for any given row. We then iterate through each column of the current row. If the column is a possible assignment of the current row, we first set `no_one` to false since the current row has at least one possible assignment.

The next step is to look through all neighbors of the subgraph vertex corresponding to the current row. Since our subgraph consists only of two vertices, the set of neighbors in this case contains a single member. We use an iterator and a for loop to iterate through the set of neighbors of the current row's vertex. This set is provided by the graph class's `neighbors` method, which returns the set of neighbors of the vertex passed to it. For each of the row's neighbors n_i , we iterate through its possible assignments. We use the boolean value `has_corresponding_neighbor` to track whether or not the row's neighbor has a possible assignment which is a neighbor of the vertex corresponding to the current column. At each possible assignment of n_i , we check whether or not the corresponding vertex in the search graph has an edge connecting it to the current column, which is in turn a possible assignment of the current row. If such an edge exists, then we mark `has_corresponding_neighbor` true, and check the next neighbor.

If the current row has no neighbors which are neighbors of the current column, then `has_corresponding_neighbor` is never set to true, so we set the current column and row of possible assignments to false, and set `changes_made` to true.

Unfortunately for us, there are no invalid assignments in the possible assignments matrix, so the refinement procedure returns true without changing the matrix. Here it is again, for convenience:

$$\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

We thus continue in `search`. Since the refinement procedure returned true, the if block is not executed, and we instead go onto a `do..while` loop. This was structured as a `do..while`, since we may need to iterate several times, but we would still like to execute the block at least once. We would later find that this was irrelevant, since the condition we used was `while(true)`, an infinite loop. Our reasoning for this should soon become apparent.

The next step is to initialize several variables for use in loops; nothing worth calling attention to. We enter a for loop to iterate through columns of the current row (the current row being given by `depth-1`), initializing the iterating variable `col` to `col_depth[depth-1]`. We stop iterating once `col < numcols`, the latter of which should be self-explanatory. Once we find a possible assignment from the current row to a column, we check if that column has

already been used by checking `cols_used[col]`. If so, we continue iterating to find another column. If not, we break out of the loop.

The first thing we do after breaking from the loop is to check whether `col` is equal to `numcols`. If so, this means we couldn't find a valid column to assign the current row to, so we must backtrack to the previous row to attempt to assign it to a different column. We do this by decrementing `depth` by 1, setting the possible assignments matrix to the one stored at `assignments_tree[depth-1]` (note that this is the new value of `depth`), and returning false. If we are at `depth 1`, then there is no previous row to reassign, and we simply return false. In our particular example, since we haven't used any columns and every entry in the possible assignments matrix is true, `col` is 0 after exiting the for loop, and `numcols` is obviously 3. Thus, we do not enter the block of the if statement.

Since we have found a valid column to assign to the current row, we set all other columns in the current row of possible assignments to false, along with all other rows in the current column. For us, that means all other columns in row 0 and all other rows in column 0. This yields the following matrix of possible assignments:

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix}$$

With one row assigned, we must recursively search for an isomorphism under this new matrix of possible assignments. We do this by calling `search` again, but with modified arguments. Namely, we append the current matrix of possible assignments to the assignments tree variable in case we need to backtrack later, we set `cols_used[col]` to true to indicate that the column `col` has been used, set `col_depth[depth-1]` to `col`, and increment `depth`. At this point, we check if `depth` is equal to p_α and, if so, return true, since that means every row has been assigned to a column. In our case, since there are only two rows, and `depth` is equal to 2, we return true. Note that this still leaves extra possible assignments in the bottom row of the possible assignments matrix. We simply ignore these and assign the bottom row to the first possible assignment in it.

If we'd had more rows, then the if block would not have executed, and we would have instead entered a different if statement which runs `search` recursively using the modified variables as arguments. If this search returns true, we have found an isomorphism under the modified possible assignments matrix, and we return true to report this finding. If it does not return true, then we set `cols_used[col]` back to false and return to the top of the while loop to try a new column. Since we had already set `col_depth[depth-1]` to `col`, we do not have to worry about the algorithm attempting to assign this row to the same column with the same result. It will continue attempting to reassign `col` starting from the next column.

The final state of our possible assignments matrix after returning is thus:

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix}$$

Since we returned true from `search`, we extract the set of possible assignments at the lowest depth in the tree that we ended on. That is, we set the assignments list to `assignments_tree[depth-2]`. We subtract 2 from `depth`, since we incremented `depth` by 1 before returning from `search`, so the list of assignments we want is actually at index `depth-2`. After assigning this to a variable we call `assignments`, we create a new vector called `isomorphism`, which contains the actual mappings which will be returned. To populate `isomorphism`, we iterate through the `assignments` matrix. Whenever we come upon an assignment, we push the current column number onto the `isomorphism` vector. In this way, `isomorphism[i]` is the vertex in the search graph assigned to vertex `i` in the subgraph. After pushing the column number, we break to go to the next row. Since each row can only have one possible assignment, there is no need to continue iterating through columns after we've found the assignment for a particular row. This also means that the bottom row, which may have multiple assignments, will have only the first of its assignments added to the vector, with the rest ignored.

Finally, after getting the assignments into a vector, we return the `isomorphism` vector. It is worth noting that, had we returned false from `search`, the if block would not have executed, and we would have instead returned an empty vector to indicate no isomorphism exists.

Back in the main function, we first check if the returned vector is empty. If so, no isomorphism was found, and we take action accordingly. In this case, we simply print a message indicating that no isomorphism was found. However, an isomorphism *was* found, so the return vector is a vector of assignments. We thus print a message indicating that an isomorphism was found, and print out which vertices in the subgraph map to which vertices in the search graph. We can retrieve the original values of the vertices given their indices by using the graph class method `get_value`, which takes the index of a vertex and returns the value of the vertex at that index in the vector of vertices in the graph. For example, for our subgraph, the vertex 10 has index 0, so `get_value(0)` will return 10. This way, we do not need to actually store the vector values along with the `isomorphism` vector; we can simply retrieve them afterwards. In our case, we have found that subgraph vertex 10 can be mapped to search graph vertex 1, and subgraph vertex 20 can be mapped to search graph vertex 2.

Since the graphs are quite small, it is trivial to draw them and verify that this result is correct. Indeed, it is one of a number possible correct results, but we only return the first one found. Future improvements to the project may include listing all possible isomorphisms, or at the very least counting them.

5 Performance

For all of our measurements, we used the worst-case input: the Facebook graph mentioned previously as the search graph, and that same graph (or, equivalently, any isomorphic graph of the same size) as the subgraph. This would take the longest of the graphs available to us, since it requires the algorithm to assign the greatest number of vertices (all of them) before returning, and, since the graphs are indeed isomorphic, there is no chance of `refine_possible_assignments` returning false, which would allow us to skip a large amount of work.

5.1 Python Implementation

The Python implementation was very inefficient, but it was necessary for us to understand what we had to do for a C++ implementation. It took between 13 and 16 seconds to determine if an isomorphism existed.

We used the Python `cProfile` library to measure the performance of our Python implementation. The documentation is available at <https://docs.python.org/3.5/library/profile.html>. A typical run of the `find_isomorphism` function timed by `cProfile` produces the following output:

```
Ordered by: standard name

ncalls  tottime  percall  cumtime  percall filename:lineno(function)
      1      0.000      0.000      15.085      15.085 <string>:1(<module>)
     492      0.000      0.000       0.001      0.000 graph.py:22(<lambda>)
  117612     13.864      0.000     13.864      0.000 graph.py:69(has_edge)
       1      0.000      0.000      15.085     15.085 isomorphism.py:43(find_isomorphism)
       1      1.218      1.218      15.085     15.085 isomorphism.py:63(update_possible_assignments)
       1      0.000      0.000      15.085     15.085 isomorphism.py:7(search)
     492      0.000      0.000       0.000      0.000 {len}
       1      0.000      0.000       0.000      0.000 {method 'disable' of '_lsprof.Profiler' objects}
     242      0.000      0.000       0.000      0.000 {method 'remove' of 'list' objects}
     488      0.002      0.000       0.002      0.000 {range}
```

Figure 5.1.3: Sample run of the Python implementation of subgraph isomorphism.

We see that, as previously explained, the overwhelming majority of the time was spent in `update_possible_assignments`, specifically running `has_edge`, which in turn calls the `in` function. The `has_edge` method is called over 117 thousand times. This is the key bottleneck in the algorithm, and any amount of improvement in that particular aspect is likely to lead to huge performance improvements, given how many times `has_edge` is called.

5.2 C++ Implementation

5.2.1 Performance Tradeoffs

While our graph class is remarkably efficient at determining whether a particular edge exists between two vertices, as well as determining the degree of a particular vertex, this efficiency is not without drawbacks. Most notably, storing

a binary matrix of edges brings the space complexity of our graph class up to $O(n^2)$. Memory is very cheap in this day and age, so sacrificing it for performance (which, as we will see, is much better) was acceptable for our purposes. This is, however, something to keep in mind for future implementations.

Additionally, updating this matrix whenever an edge or vertex is added to the graph adds considerable overhead in itself, making `add_edge` and `add_vertex` more expensive methods to run.

5.2.2 Results

The source code used to test the performance of the C++ implementations can be found in the appendix. We used Linux's `perf_event_open` system call, which allows for high-precision performance monitoring. The only part of the program which is timed is the call to `find_isomorphism`.

The C++ implementation, due to the improvement of `has_edge`'s runtime, runs significantly more quickly than the Python implementation; it runs in about 8 seconds, compared to about 16 seconds for the Python implementation. This is about twice as fast, even without any parallelization. However, we found it was generally more precise to measure the number of CPU clock cycles taken for the algorithm to complete, so we measured that instead. The output of the program used to measure this is below:

```
Isomorphism found! Assignments omitted, since we're just checking performance.
Searching for isomorphism took 26142295153 CPU cycles
```

Figure 5.2.1: Performance of non-parallelized version of our algorithm.

We see that the C++ version took 26,142,295,153 cycles to complete on the system used to test. This is equivalent to about eight seconds, as mentioned previously. The exact number of cycles will vary slightly from trial to trial, but repeated trials have made us confident that this is a representative sample.

5.3 OpenMP Optimization

The OpenMP code performs very well compared to the non-parallelized version. The system we tested it on has a four-core CPU, allowing for up to four threads to run in parallel, leading to a four-fold performance improvement. It is important to note, however, that even in the case of absolute perfect parallelization, it is impossible to achieve exactly four times the performance. Using OpenMP requires the program execution to split up into four threads and combine back into one after each of the four threads have finished the work assigned to them. This necessitates the execution of a few additional instructions beyond the ones relevant to the computation in our algorithm. Thus, even in the best case, we can only hope for performance which is slightly less than four times faster.

Moreover, we must consider the fact that we are not parallelizing the entire algorithm; only `refine_possible_assignments` is run in parallel, and only part of it, at that. The rest of the algorithm runs serially, like the non-parallelized version. Thus, even without the overhead involved with using OpenMP, we cannot ever expect the entire algorithm to run four times faster, when only part of it runs (theoretically) four times as fast.

With all that said, we came surprisingly close to four times faster, all things considered. The output of the program used to measure the performance of the OpenMP code is below:

```
Isomorphism found! Assignments omitted, since we're just checking performance.
Searching for isomorphism took 7125381701 CPU cycles
```

Figure 5.3.1: Performance of OpenMP version of our algorithm.

Here we see instead that the algorithm now takes only 7,125,381,701 cycles to complete. Once again, this number will vary slightly for each run, but this particular number was about average. Using the number from figure 5.2.1, we can compute how much faster this particular run of the parallelized algorithm is than that particular run of the non-parallelized algorithm by dividing the latter by the former. This comes out to:

$$\frac{26142295153\text{cycles}}{7125381701\text{ cycles}} \approx 3.669$$

At 3.669 times faster, the parallel version of the algorithm runs much closer to four times faster than we had anticipated, given the factors discussed earlier. This greatly supports our hypothesis that parallelizing just the refinement procedure leads to a much more efficient implementation, since that function is by far the most computationally expensive portion of the algorithm, and is called many, many times.

6 Conclusion

Using parallelization techniques in order to increase efficiency of the NP-complete problem of subgraph isomorphism proved to be a success. In our C++ implementation of Ullman's algorithm with OpenMP, graph isomorphism took two to three seconds to compute; this is down from the average of eight seconds from the C++ implementation without parallelization, using what was easily our worst-case scenario. Additionally, solving smaller instances of the subgraph isomorphism problem take less than a second. In either case, using C++ and adding OpenMP optimizations resulted in a much faster implementation than our Python implementation, which took between 14 and 16 seconds to run.

When it comes to innovation, everyone wants their next thing to be more efficient, both time-wise and space-wise.

This project has shown that even the smallest amount of parallelization can make a large difference for very complex problems. If we had combined various parallelization techniques, such as utilizing SIMD instructions alongside the multithreading technique we used, we might have gotten an implementation even faster than what we had. Parallelization is seeing growing usage to address the problems of an increasingly demanding world.

7 Appendix

7.1 Python Code

7.1.1 Graph Class

```
"""
Graph module
only has a class
"""

class Vertex:
    """
    Vertex Class
    """
    def __init__(self, label = "", data="", visited=False):
        self.label = label
        self.data = data
        self.visited = visited

class Graph:
    """
    Graph class
    """
    def __init__(self, vertices=None, edges=None, adjacencies=None):
        if vertices is None:
            self.vertices = []
            self.n_vertices = lambda : len(self.vertices)
        else:
```

```
        self.vertices = vertices
    if edges is None:
        self.edges = []
    else:
        self.edges = edges
    if adjacencies is None:
        self.adjacencies = []
    else:
        self.adjacencies = adjacencies

    self.__update_adjacencies()

#Private, don't want anyone calling this by accident.
def __update_adjacencies(self):
    for edge in self.edges:
        self.adjacencies.append({edge[0], edge[1]})

def add_vertex(self, vertex=Vertex()):
    """ Add vertex to graph
    """
    self.vertices.append(vertex)

def remove_vertex(self, vertex):
    """ Remove vertex from graph
        along with edges which contain it
    """
    self.vertices.remove(vertex)
    # list comprehension
    # TLDR remove all elements in self.edges which contain vertex
    # we need this since edges is a list of tuples, not of single elements
    self.edges.remove([edge for edge in self.edges if vertex in edge])
    # see above
```

```
self.adjacencies.remove([adj for adj in self.adjacencies if vertex in adj])

def add_edge(self, source, dest, weight):
    """ Add edge connecting source and dest
    along with weight
    """
    if (source or dest) not in self.vertices:
        self.vertices.append(source)
        self.vertices.append(dest)

    self.edges.append([source, dest, weight])
    self.adjacencies.append({source, dest})

def has_edge(self, vert1, vert2):
    """ Checks if edge connecting vert1 and vert2 is in the graph
    """
    #if adjacent, there's an edge
    return ({vert1, vert2} in self.adjacencies)

# for debugging
def print_graph(self):
    """ Prints information about the graph
    """
    print("Graph (vertices, edges, adjacencies):")
    print(self.vertices)
    print(self.edges)
    print(self.adjacencies)
```

Listing 1: Python graph class.

7.1.2 Subgraph Isomorphism Algorithm

```
""" Subgraph Isomorphism
```


CPU Implementation

"""

```
import copy as cp # for deepcopy
```

```
def search(graph, subgraph, assignments, possible_assignments):
```

```
    """ Uses DFS to find instance of subgraph within larger graph
    """
```

```
    update_possible_assignments(graph, subgraph, possible_assignments)
```

```
    i = len(assignments)
```

```
    # Make sure that every edge between assigned vertices in the subgraph
    # is also an edge in the graph.
```

```
    # this loop calls graph.has_edge once for every edge in the subgraph
```

```
    for edge in subgraph.edges:
```

```
        if edge[0] < i and edge[1] < i:
```

```
            if not graph.has_edge(assignments[edge[0]], assignments[edge[1]]):
                return False
```

```
    # If all the vertices in the subgraph are assigned, then we are done.
```

```
    if i == subgraph.n_vertices():
```

```
        return True
```

```
    for j in possible_assignments[i]:
```

```
        if j not in assignments:
```

```
            assignments.append(j)
```

```
    # Create a new set of possible assignments, where graph node j
```

```
    # is the only possibility for the assignment of subgraph node i.
```



```
        # for every item in the subgraph's adjacencies
        # for every possible assignments
        # that is a huge number of calls to has_edge
        # which is in itself an O(n) operation
        # definitely room for improvement.
        if adj in possible_assignments[adj] and \
            graph.has_edge(j, vert):
            match = True
    if not match:
        possible_assignments[i].remove(j)
        any_changes = True
```

Listing 2: Python algorithm to find isomorphism.

7.1.3 Main

```
""" Main
"""

import sys #args
import getopt
from graph import Graph
from isomorphism import find_isomorphism
import cProfile

def main(argv):
    """ Main function
    """
    argc = len(argv)

    interactive, debug = handle_args(argc, argv)

    if interactive:
        print("Enter the graph as a space-separated adjacency list of the form\
```

```
        node1 node2 [edge_weight]:")
    filename = sys.stdin
else:
    filename = argv[1]

graph = import_data(filename, debug)

if debug:
    graph.print_graph()

if interactive:
    print("Enter the subgraph using the same format:")
    sub = sys.stdin
else:
    sub = argv[2]

subgraph = import_data(sub, debug)

if debug:
    subgraph.print_graph()

cProfile.runctx('find_isomorphism(graph, subgraph)', globals(), locals())

def import_data(filename, debug):
    """ Takes graph adjacency list written in file fd
        creates and returns graph object representing
        the graph in the file
        Note that edge weights are all 0 unless otherwise specified.
    """
    if filename == sys.stdin:
        fdesc = filename
    else:
```

```
    try:
        fdesc = open(filename, 'r')
    except IOError as err:
        sys.stderr.write("Error: Could not open file '%s'. Aborting.\n%s\n" \
                        % (filename, err))
        sys.exit(2)

if debug:
    print(fdesc)

graph = Graph()

if debug:
    print("Reading file")
for line in fdesc:
    vals = line.split()
    if len(vals) == 3:
        edge_weight = int(vals[2])
    elif len(vals) > 2:
        sys.stderr.write("Error reading file '%s'. Data should be in form \
                        'node1 node2 [edge_weight]', but line was %s\n" % line)
        sys.exit(3)
    else:
        edge_weight = 0

    graph.add_edge(int(vals[0]), int(vals[1]), edge_weight)

if filename != sys.stdin:
    fdesc.close()

return graph
```

```
def usage(name, err):
    """ Prints usage info
    """
    usg = """usage: %s [FILE] [OPTION]
Options and arguments:
FILE:                Space-separated file containing graph adjacency list. \
                        If edge weights are not specified, all weights will be assumed 0.
-h, --help:          Print this help message and exit.
-i, --interactive:    Get space-separated graph adjacency list from stdin.
-d, --debug:          Print debugging information.\n"""
    if err:
        sys.stderr.write(usg % name)
        return

    sys.stdout.write(usg % name)

def handle_args(argc, argv):
    """ get command line arguments, set variables accordingly
    """
    # TODO change to False for final deployment
    debug = True
    interactive = False

    # arguments should be either 2 files or the two files and some args
    # unless -i is specified in which case just the args
    if argc < 2:
        usage(argv[0], True)
        sys.exit(1)
    else:
        try:
            # using _ as a variable name is sort of the conventional way of saying
            # "we don't need or use this variable but i'm assigning it because
```

```
# the method or api or whatever i'm using requires it"
_, args = getopt.getopt(argv, "hid", ["help", "interactive", "debug"])
except getopt.GetoptError:
    # print usage info and quit
    usage(argv[0], True)
    sys.exit(1)
# handle args
for arg in args:
    if arg in ("-h", "--help"):
        usage(argv[0], False)
        sys.exit()
    elif arg in ("-i", "--interactive"):
        interactive = True
        args.remove(arg)
    elif arg in ("-d", "--debug"):
        debug = True
        args.remove(arg)

return interactive, debug

if __name__ == "__main__":
    main(sys.argv)
```

Listing 3: Python main function to process arguments and run the isomorphism algorithm.

7.2 C++ Code

7.2.1 Graph Class

```
#include <algorithm> //Replace, remove, find
#include <iostream>
#include "graph.h"

using std::remove;
```

```
using std::find;
using std::get;
using std::cout;
using std::map;
using std::unordered_set;

struct CompareFirst
{
    CompareFirst(int val) : val_(val) {}

    bool operator()(const std::pair<int, int>& elem) const {
        return val_ == elem.first;
    }

private:
    int val_;
};

// edges[i][j] = weight of edge connecting i and j
// -1 if no such edge

Graph::Graph () {
    this->vertices = vector<pair<int, int> >();
    this->edges = vector<vector<int> >();
    this->adjacencies = vector<unordered_set<int> >();
    this->vertex_indices = map<int, int>();
    this->vertex_vals = map<int, int>();
}

void Graph::add_vertex(pair<int, int> vertex) {
    if (this->get_index(vertex.first) < 0) {
        this->vertices.push_back(vertex);
        // add mappings
        //vertex_indices[vertex.first] = vertices.size()-1;
    }
}
```



```
        this->vertex_indices.insert(pair<int, int>(vertex.first, vertices.size()-1));

        //vertex_vals[vertices.size()-1] = vertex.first;
        this->vertex_vals.insert(pair<int, int>(vertices.size()-1, vertex.first));

        // add row with edges.size()+1 -1s
        this->edges.push_back(vector<int>(edges.size()+1,-1));
        // add a set for its neighbors
        this->adjacencies.push_back(unordered_set<int>());
    }
}

void Graph::add_edge(int source, int dest, int weight) {

    // WEIGHT SHOULD BE NON-NEGATIVE
    this->add_vertex(pair<int, int>(source, 0));

    this->add_vertex(pair<int, int>(dest, 0));

    int sourceind = this->get_index(source);
    int destind = this->get_index(dest);

    // increase vertex degree
    this->vertices[this->get_index(source)].second++;
    this->vertices[this->get_index(dest)].second++;

    // source -> dest
    if (sourceind < destind) {
        this->edges[destind][sourceind] = weight;
    }

    else {
```

```
        this->edges[sourceind][destind] = weight;
    }

    // TODO: get adjacency set of source and dest (or create them if not in graph)
    // add dest/source to that set
    // once again, consider directed graphs
    this->add_neighbor(source, dest);
    this->add_neighbor(dest, source);
}

bool Graph::has_edge(int vert1, int vert2){
    // just got this bad boy down to O(1)
    // you're welcome
    return (edges[vert1][vert2] >= 0);
}

int Graph::get_index(int value) {
    if (this->vertex_indices.find(value) != this->vertex_indices.end()) {
        return this->vertex_indices[value];
    }

    return -1;
}

int Graph::get_value(int index) {
    if (this->vertex_vals.find(index) != this->vertex_vals.end()) {
        return this->vertex_vals[index];
    }

    return -1;
}
```

```

unordered_set<int> Graph::neighbors(int vertex) {
    return this->adjacencies[this->get_index(vertex)];
}

void Graph::add_neighbor (int vert, int neighbor) {
    this->neighbors(vert).insert(this->vertices[this->get_index(neighbor)].first);
}

// basically obsolete now
//void Graph::update_adjacencies() {
//    for (size_t i = 0; i < edges.size(); i++){
//        unordered_set<int> edge ({get<0>(edges[i]), get<1>(edges[i])});
//        adjacencies.push_back(edge);
//    }
//}

// void Graph::print_graph() {
//    cout << "Begin graph" << std::endl;
//    cout << "-----\n";
//    cout << "Vertices: " << vertices << std::endl;
//    cout << "Edges: " << edges << std::endl;
//    cout << "-----\n";
//    cout << "End graph" << std::endl;
// }

```

Listing 4: C++ graph class.

7.2.2 Subgraph Isomorphism Algorithm

```

#include <vector>
#include <utility>
#include <unordered_set>
#include <algorithm> //for std::find

```

```

#include "isomorphism.h"

using std::find;

vector<int> find_isomorphism (Graph &sub, Graph &graph) {

    vector< vector<bool> > > possible_assignments = create_possible_assignments(sub, graph);
    vector<vector< vector< bool> > > > assignments_tree;

    // we need a way to keep track of which columns we've already assigned for each row
    // cols_used[i] == true if column i has already been used
    vector<bool> cols_used = vector<bool>(possible_assignments[0].size(), false);
    // keep track of which column was used at which depth
    // col_depth[i-1] == k iff column k was used at depth i
    vector<size_t> col_depth = vector<size_t>(possible_assignments.size(), 0);

    size_t depth = 1;

    if (search(sub, graph, assignments_tree, possible_assignments, cols_used, col_depth, depth)) {

        vector<vector<bool> > > assignments = assignments_tree[depth-2];
        vector<int> isomorphism;
        for (size_t i = 0; i < assignments.size(); i++) {
            //fprintf(stderr, "Row %lu\n", i);
            for (size_t j = 0; j < assignments[i].size(); j++) {
                //fprintf(stderr, "%d\n", assignments[i][j] ? 1 : 0);
                if (assignments[i][j]) {
                    isomorphism.push_back(j);
                    // go to next row
                    // no need to waste iterations
                    break;
                }
            }
        }
    }
}

```

```
    }

    return isomorphism;
}

return vector<int>(0);
}

bool search (Graph &sub, Graph &graph,
             vector<vector<vector<bool> > > &assignments_tree,
             vector<vector<bool> > &possible_assignments,
             vector<bool> &cols_used,
             vector<size_t> &col_depth,
             size_t &depth) {

    // if refine fails, no possible isomorphism
    if (!refine_possible_assignments(sub, graph, possible_assignments)) {
        return false;
    }

    do {
        size_t numRows = possible_assignments.size(), numcols = possible_assignments[0].size(),

        for (col = col_depth[depth-1]; col < numcols; col++) {
            // find first column we havent already used
            // col is index of vertex in search graph we're trying to assign
            // we start searching from col_depth[depth-1]+1 since
            // that's the col after the one we attempted previously and
            // we gotta find the next one
            if (possible_assignments[depth-1][col] && !cols_used[col]) {
                break;
            }
        }
    }
```

```
}

if (col == numcols) {
    // we couldn't find a valid column.
    // if we're on depth 1, there's no level to go back up to
    // no possible isomorphism
    if (depth == 1) {
        return false;
    }

    // otherwise, go up to the previous depth and try to assign that
    // row to another column
    else {
        depth--;
        possible_assignments = assignments_tree[depth-1];
        return false;
    }
}

// we found a valid column
// set all other possible assignments in this col to 0
// if we're assigning this vertex to the subgraph vertex
// we can't have any other subgraph vertex also assigned to it
size_t i;
for (i = 0; i < numrows; i++) {
    if (i != depth-1) {
        possible_assignments[i][col] = false;
    }
}

// set all other 1s in this row to 0
// we can only assign each subgraph vertex to a single column
```

```
for (i = 0; i < numcols; i++) {
    if (i != col) {
        possible_assignments[depth-1][i] = false;
    }
}

// now go to the next row and repeat
// we used this column
cols_used[col] = true;
// and we used it at this depth
col_depth[depth-1] = col;
// save a copy of possible assignments
assignments_tree.push_back(possible_assignments);
// increment depth
depth++;

if (depth == possible_assignments.size()) {
    // we've assigned every row to a column
    // isomorphism found
    return true;
}

// if we get all the way to the bottom by doing this repeatedly,
// we'll have found an isomorphism
if (search(sub, graph, assignments_tree, possible_assignments, cols_used, col_depth, depth)) {
    return true;
}

else {
    // that didn't work. find a new column.
    cols_used[col] = false;
}

// this is fine since we only get here if we need to find a new column
```

```

    } while (true);
}

vector < vector<bool> > create_possible_assignments(Graph &sub, Graph &graph) {
    // possible_assignments[i][j] == true iff a possible assignment exists from i in sub
    // to j in search graph
    vector < vector<bool> > possible_assignments (sub.vertices.size(),
        vector<bool>(graph.vertices.size(), false));

    // at first, every vertex in search graph with rank >= i is a possible assignment
    // this will be refined later
    // see refine_possible_assignments
    size_t i, j;
    size_t s_n = sub.vertices.size();
    size_t g_n = graph.vertices.size();
    for (i = 0; i < s_n; i++) {
        for (j = 0; j < g_n; j++) {
            if (graph.vertices[j].second >= sub.vertices[i].second) {
                possible_assignments[i][j] = true;
            }
        }
    }

    return possible_assignments;
}

bool refine_possible_assignments(Graph &sub, Graph &graph, vector < vector<bool> > &possible_assignments) {
    size_t i, j;
    size_t pa_n = possible_assignments.size();
    size_t pb_n = possible_assignments[0].size();
    bool changes_made = true;

```



```
while (changes_made) {
    changes_made = false;
    for (i = 0; i < pa_n; i++) {
        // check if this row contains no 1s
        bool no_one = true;
        for (j = 0; j < pb_n; j++) {
            if (possible_assignments[i][j]) {
                no_one = false;

                // check if all of i's neighbors have a possible assignment to a neighbor of j
                // iterate through all neighbors of i
                unordered_set<int> neighbors_i = sub.neighbors(sub.get_value(i));
                unordered_set<int>::iterator n_i = neighbors_i.begin();

                for (; n_i != neighbors_i.end(); n_i++) {
                    // for each neighbor of i, iterate through all its possible assignments
                    size_t k;
                    bool has_corresponding_neighbor = false;
                    for (k = 0; k < pb_n; k++) {
                        if (possible_assignments[sub.get_index(*n_i)][k]) {
                            // for each possible assignment from i's neighbor to graph, check if it is
                            if (graph.has_edge(graph.get_value(j), graph.get_value(k))) {
                                // if so, then check the next neighbor
                                has_corresponding_neighbor = true;
                                break;
                            }
                        }
                    }
                }

                // if this neighbor has no corresponding neighbor, then j is an invalid match.
                // move on to the next possible assignment
                if (!has_corresponding_neighbor) {
                    possible_assignments[i][j] = 0;
                }
            }
        }
    }
}
```

```
        changes_made = true;
        break;
    }
}
}

// we never found a 1 in this row, so there's no point in continuing
// there's at least one vertex in subgraph that cannot be mapped to any
// vertex in the search graph, so we know that
// possible_assignments cannot specify any isomorphism
if (no_one) {
    return false;
}

}

// M was successfully refined without creating any rows with all 0s. return true.
return true;
}
```

Listing 5: C++ algorithm to find isomorphism.

7.2.3 Main

```
// use stderr to print debug info since stderr is unbuffered
// i.e. anything sent to stderr will be printed immediately
// instead of put into a buffer and printed when the buffer
// gets flushed
#include <stdio>
#include <stdlib>
#include <cstring>
#include <iostream> //debugging
```

```
#include <vector>
#include <tuple>
#include <getopt.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <unistd.h>
#include <fcntl.h>
#include <errno.h>
#include "graph.h"
#include "isomorphism.h"
#include <unordered_set>
#include <cstring>
using std::string;

const char* program_name;

vector<string> handle_args(int argc, char **argv);
void usage (FILE* stream, int exit_code);
Graph import_data(const char *filename, const int debug);

int main(int argc, char **argv) {

    program_name = argv[0];

    vector<string> args = handle_args(argc, argv);

    int interactive = atoi(args[0].c_str());
    int debug = atoi(args[1].c_str());

    char filename_default[] = "stdin";

    const char *graph_filename = (!interactive) ? args[2].c_str()
```

```
                                : filename_default;

const char *sub_filename = (!interactive) ? args[3].c_str()
                                : filename_default;

if (debug) {
    fprintf(stderr, "inter = %d, debug = %d\n", interactive, debug);
    fprintf(stderr, "Graph file %s, subgraph file %s\n", graph_filename, sub_filename);
}

Graph graph = import_data(graph_filename, debug);

if (debug) {
    fprintf(stderr, "Successfully imported graph from file %s.\n",
                graph_filename);
}

Graph subgraph = import_data(sub_filename, debug);

if (debug) {
    fprintf(stderr, "Successfully imported subgraph from file %s.\n",
                sub_filename);
}

if (debug) {
    fprintf(stderr, "Running find_isomorphism....\n");
}

vector<int> isomorphism = find_isomorphism(subgraph, graph);

if (debug){
    fprintf(stderr, "find_isomorphism terminated.\n");
}
```

```
    }

    if (isomorphism.size() == 0) {
        printf("No isomorphism found.\n");
    }

    else{
        printf("Isomorphism found! Assignments are:\n");
        size_t num_assignments = isomorphism.size();
        for (size_t i = 0; i < num_assignments; i++) {
            int subgraph_vert = subgraph.get_value(i);
            int graph_vert = graph.get_value(isomorphism[i]);
            printf("Vertex %d maps to vertex %d\n", subgraph_vert, graph_vert);
        }
    }

    return 0;
}
```

```
Graph import_data(const char *filename, const int debug) {
```

```
    FILE *fd;

    if (strcmp(filename, "stdin") == 0) {
        fd = stdin;
    }

    else {
        fd = fopen(filename, "r");
        if (fd == NULL) {
            fprintf(stderr, "Could not open file %s for reading.\n", filename);
            perror("Error calling fopen()");
            exit(errno);
        }
    }
}
```

```
    }
}

char *line = NULL;
size_t n = 0;

Graph g;

if (debug) {
    fprintf(stderr, "Reading file %s...\n", filename);
}

char *tempa, *tempb, *tempweight;
char tempweight_default[] = "0";
ssize_t bytes_read;
// getline returns -1 on failure to read a line (including EOF)
bytes_read = getline(&line, &n, fd);
while (bytes_read != -1) {
    tempa = strtok(line, " ");
    // strtok returns NULL if the token (space in this case) was not found in the str
    // otherwise return the string up to the token (exclusive)
    if (tempa == NULL) {
        fprintf(stderr, "Error parsing file %s. Adjacency list must be of form 'node1 node2\n"
            "but line was %s\n", filename, line);
        exit(-1);
    }

    // if arg is null after a call to strtok, it'll keep reading after the position of the
    tempb = strtok(NULL, " ");
    if (tempb == NULL) {
        fprintf(stderr, "Error parsing file %s. Adjacency list must be of form 'node1 node2\n"
            "but line was %s\n", filename, line);
    }
}
```

```
        exit(-1);
    }

    tempweight = strtok(NULL, " ");
    if (tempweight == NULL) {
        tempweight = tempweight_default;
    }

    int atempa = atoi(tempa);
    int atempb = atoi(tempb);
    int atempweight = atoi(tempweight);

    // this will also add the vertices if they don't exist
    // add_vertex is called from within add_edge if no vertex
    // with the given value exists
    g.add_edge(atempa, atempb, atempweight);
    bytes_read = getline(&line, &n, fd);
}

if (fd != stdin) {
    int stat = fclose(fd);
    if (stat != 0) {
        perror("Error closing file");

        // we never wrote to the file so we probably don't even need to
        // check for errors in the first place
        // but just kill the program if there's an error closing it
        exit(errno);
    }
}

// don't free line since it's null now
```

```
    return g;
}

vector<string> handle_args(int argc, char **argv) {
    // for getopt
    const char* const short_options = "hid";
    const struct option long_options[] = {
        { "help", 0, NULL, 'h' },
        { "interactive", 0, NULL, 'i' },
        { "debug", 0, NULL, 'd' },
        { NULL, 0, NULL, 0 }
    };

    vector<string> returnargs;
    returnargs.resize(4);
    int interactive = 0;
    int debug = 0;

    if (argc < 2) {
        usage(stderr, 1);
    }

    int next_option;
    do {
        next_option = getopt_long(argc, argv, short_options,
                                long_options, NULL);

        switch(next_option) {
            case 'h':
                usage(stdout, 0);
                // calling usage quits the program anyway but whatever
                // we'll put a break here anyway
                break;
        }
    } while (next_option != -1);

    return returnargs;
}
```



```
    case 'i':
        interactive = 1;
        break;

    case 'd':
        debug = 1;
        break;

    case '?': // invalid option
        usage(stderr, 1);
        break;

    case -1: // no more options
        break;

    default: // something went very wrong if we got here
        abort();
}
}
while (next_option != -1);

// now that we've gone through all the options, OPIND points to first
// nonoption argument. which is hopefully the first filename. unless
// interactive was specified, in which case yell at the user.

if (!interactive) {
    returnargs[2] = string(argv[optind]);
    returnargs[3] = string(argv[optind+1]);
}
```

```

if (debug) {
    fprintf(stderr, "Interactive set to %d. Debug is %d.\n", interactive, debug);
    if (!interactive) {
        fprintf(stderr, "Using file %s for graph and %s for subgraph.\n",
            argv[optind], argv[optind+1]);
    }
}

// put what we want to return in their places.
returnargs[0] = std::to_string(interactive);
returnargs[1] = std::to_string(debug);

return returnargs;
}

void usage (FILE* stream, int exit_code) {
    fprintf(stream, "usage: %s [FILE] [OPTION]\n", program_name);
    fprintf(stream,
        "Options and arguments:\n"
        "FILE:           Space-separated file containing graph adjacency list. If edge weights are present, they must be separated by a space.\n"
        "-h, --help:      Print this help message and exit.\n"
        "-i, --interactive: Get space-separated graph adjacency list from stdin.\n"
        "-d, --debug:      Print debugging information.\n");

    exit(exit_code);
}

```

Listing 6: C++ main function to process arguments and run the isomorphism algorithm.

7.3 OpenMP Code

Note: Since the OpenMP version of our algorithm uses the same graph class and `main.cpp`, we have omitted them from this section.

7.3.1 Subgraph Isomorphism Algorithm

```

#include <vector>
#include <utility>
#include <unordered_set>
#include <algorithm> //for std::find
#include "isomorphism.h"
#include "omp.h"

int THREADS = omp_get_max_threads();

using std::find;

vector<int> find_isomorphism (Graph &sub, Graph &graph) {

    vector< vector<bool> > > possible_assignments = create_possible_assignments(sub, graph);
    vector<vector< vector< bool> > > > assignments_tree;

    // we need a way to keep track of which columns we've already assigned for each row
    // cols_used[i] == true if column i has already been used
    vector<bool> cols_used = vector<bool>(possible_assignments[0].size(), false);
    // keep track of which column was used at which depth
    // col_depth[i-1] == k iff column k was used at depth i
    vector<size_t> col_depth = vector<size_t>(possible_assignments.size(), 0);

    size_t depth = 1;

    if (search(sub, graph, assignments_tree, possible_assignments, cols_used, col_depth, depth)) {

        vector<vector<bool> > > assignments = assignments_tree[depth-2];
        vector<int> isomorphism;

        for (size_t i = 0; i < assignments.size(); i++) {
            //fprintf(stderr, "Row %lu\n", i);
            for (size_t j = 0; j < assignments[i].size(); j++) {

```

```
        //fprintf(stderr, "%d\n", assignments[i][j] ? 1 : 0);
        if (assignments[i][j]) {
            isomorphism.push_back(j);
            // go to next row
            // no need to waste iterations
            break;
        }
    }

    return isomorphism;
}

return vector<int>(0);
}

bool search (Graph &sub, Graph &graph,
             vector<vector<vector<bool> > > &assignments_tree,
             vector<vector<bool> > &possible_assignments,
             vector<bool> &cols_used,
             vector<size_t> &col_depth,
             size_t &depth) {

    // if refine fails, no possible isomorphism
    if (!refine_possible_assignments(sub, graph, possible_assignments)) {
        return false;
    }

    do {
        size_t numRows = possible_assignments.size(), numcols = possible_assignments[0].size(),

        for (col = col_depth[depth-1]; col < numcols; col++) {
```

```
// find first column we havent already used
// col is index of vertex in search graph we're trying to assign
// we start searching from col_depth[depth-1]+1 since
// that's the col after the one we attempted previously and
// we gotta find the next one
if (possible_assignments[depth-1][col] && !cols_used[col]) {
    break;
}

if (col == numcols) {

    // we couldn't find a valid column.
    // if we're on depth 1, there's no level to go back up to
    // no possible isomorphism
    if (depth == 1) {
        return false;
    }

    // otherwise, go up to the previous depth and try to assign that
    // row to another column
    else {
        depth--;
        possible_assignments = assignments_tree[depth-1];
        return false;
    }
}

// we found a valid column
// set all other possible assignments in this col to 0
// if we're assigning this vertex to the subgraph vertex
// we can't have any other subgraph vertex also assigned to it
size_t i;
```

```
for (i = 0; i < numrows; i++) {
    if (i != depth-1) {
        possible_assignments[i][col] = false;
    }
}

// set all other 1s in this row to 0
// we can only assign each subgraph vertex to a single column
for (i = 0; i < numcols; i++) {
    if (i != col) {
        possible_assignments[depth-1][i] = false;
    }
}

// now go to the next row and repeat
// we used this column
cols_used[col] = true;
// and we used it at this depth
col_depth[depth-1] = col;
// save a copy of possible assignments
assignments_tree.push_back(possible_assignments);
// increment depth
depth++;

if (depth == possible_assignments.size()) {
    // we've assigned every row to a column
    // isomorphism found
    return true;
}

// if we get all the way to the bottom by doing this repeatedly,
// we'll have found an isomorphism
```

```

    if (search(sub, graph, assignments_tree, possible_assignments, cols_used, col_depth, de
        return true;
    }
    else {
        // that didn't work. find a new column.
        cols_used[col] = false;
    }
    // this is fine since we only get here if we need to find a new column
} while (true);
}

vector < vector<bool> > create_possible_assignments(Graph &sub, Graph &graph) {
    // possible_assignments[i][j] == true iff a possible assignment exists from i in sub
    // to j in search graph
    vector < vector<bool> > possible_assignments (sub.vertices.size(),
        vector<bool>(graph.vertices.size(), false));

    // at first, every vertex in search graph with rank >= i is a possible assignment
    // this will be refined later
    // see refine_possible_assignments
    size_t i, j;
    size_t s_n = sub.vertices.size();
    size_t g_n = graph.vertices.size();
    for (i = 0; i < s_n; i++) {
        for (j = 0; j < g_n; j++) {
            if (graph.vertices[j].second >= sub.vertices[i].second) {
                possible_assignments[i][j] = true;
            }
        }
    }

    return possible_assignments;
}

```

```

}
```

```

bool refine_possible_assignments(Graph &sub, Graph &graph, vector < vector<bool> > &possible_assignments) {
    size_t i, j;
    size_t pa_n = possible_assignments.size();
    size_t pb_n = possible_assignments[0].size();
    bool changes_made = true;

    while (changes_made) {
        changes_made = false;
        bool no_one = true;
        for (i = 0; i < pa_n; i++) {
            // check if this row contains no 1s
            no_one = true;
#pragma omp parallel num_threads(THREADS)
            for (j = 0; j < pb_n; j++) {
                if (possible_assignments[i][j]) {
                    no_one = false;
                    // check if all of i's neighbors have a possible assignment to a neighbor of j
                    // iterate through all neighbors of i
                    unordered_set<int> neighbors_i = sub.neighbors(sub.get_value(i));
                    unordered_set<int>::iterator n_i = neighbors_i.begin();

                    for (; n_i != neighbors_i.end(); n_i++) {
                        // for each neighbor of i, iterate through all its possible assignments
                        size_t k;
                        bool has_corresponding_neighbor = false;
                        for (k = 0; k < pb_n; k++) {
                            if (possible_assignments[sub.get_index(*n_i)][k]) {
                                // for each possible assignment from i's neighbor to graph, check if it is
                                if (graph.has_edge(graph.get_value(j), graph.get_value(k))) {
                                    // if so, then check the next neighbor
                                }
                            }
                        }
                    }
                }
            }
        }
    }
}
```



```
        has_corresponding_neighbor = true;
        break;
    }
}

// if this neighbor has no corresponding neighbor, then j is an invalid match.
// move on to the next possible assignment
if (!has_corresponding_neighbor) {
    possible_assignments[i][j] = 0;
    changes_made = true;
    break;
}

}

}

// we never found a 1 in this row, so there's no point in continuing
// there's at least one vertex in subgraph that cannot be mapped to any
// vertex in the search graph, so we know that
// possible_assignments cannot specify any isomorphism
if (no_one) {
    return false;
}

}

// M was successfully refined without creating any rows with all 0s. return true.
return true;
}
```

Listing 7: C++ algorithm to find isomorphism, with parallelisms through OpenMP.

7.4 Code Used to Check Performance of C++ code

```
#include <cstdio>
#include <cstdlib>
#include <cstring>
#include <getopt.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <unistd.h>
#include <fcntl.h>
#include <errno.h>
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <string.h>
#include <sys/ioctl.h>
#include <linux/perf_event.h>
#include <asm/unistd.h>
#include <cstring>
#include <vector>
#include <string>

#include "isomorphism.h"

const char* program_name;

using std::string;
using std::vector;

void set_up_graphs(int argc, char **argv, Graph &sub, Graph &graph);
vector<string> handle_args(int argc, char **argv);
void usage (FILE* stream, int exit_code);
Graph import_data(const char *filename);
```

```
static long perf_event_open(struct perf_event_attr *hw_event, pid_t pid,
    int cpu, int group_fd, unsigned long flags) {

    int ret;

    ret = syscall(__NR_perf_event_open, hw_event, pid, cpu,
        group_fd, flags);
    return ret;
}

int main(int argc, char **argv) {

    struct perf_event_attr pe;
    long long count;
    int fd;

    memset(&pe, 0, sizeof(struct perf_event_attr));
    pe.type = PERF_TYPE_HARDWARE;
    pe.size = sizeof(struct perf_event_attr);
    pe.config = PERF_COUNT_HW_CPU_CYCLES;
    pe.disabled = 1;
    pe.exclude_kernel = 1;
    pe.exclude_hv = 1;

    fd = perf_event_open(&pe, 0, -1, -1, 0);
    if (fd == -1) {
        fprintf(stderr, "Error opening leader %llx\n", pe.config);
        exit(EXIT_FAILURE);
    }

    Graph subgraph, graph;
```

```
set_up_graphs(argc, argv, subgraph, graph);

ioctl(fd, PERF_EVENT_IOC_RESET, 0);
ioctl(fd, PERF_EVENT_IOC_ENABLE, 0);

vector<int> isomorphism = find_isomorphism(subgraph, graph);

ioctl(fd, PERF_EVENT_IOC_DISABLE, 0);
read(fd, &count, sizeof(long long));

if (isomorphism.size() == 0) {
    printf("No isomorphism found.\n");
}

else {
    printf("Isomorphism found! Assignments omitted, since we're just checking"
           " performance.\n");
}

printf("Searching for isomorphism took %lld CPU cycles\n", count);

close(fd);
}

void set_up_graphs(int argc, char **argv, Graph &sub, Graph &graph) {

    program_name = argv[0];

    vector<string> args = handle_args(argc, argv);

    int interactive = atoi(args[0].c_str());
```

```
char filename_default[] = "stdin";

const char *graph_filename = (!interactive) ? args[2].c_str()
                                : filename_default;
const char *sub_filename = (!interactive) ? args[3].c_str()
                                : filename_default;

graph = import_data(graph_filename);

sub = import_data(sub_filename);

return;
}

vector<string> handle_args(int argc, char **argv) {
    // for getopt
    const char* const short_options = "hid";
    const struct option long_options[] = {
        { "help", 0, NULL, 'h' },
        { "interactive", 0, NULL, 'i' },
        { "debug", 0, NULL, 'd' },
        { NULL, 0, NULL, 0 }
    };

    vector<string> returnargs;
    returnargs.resize(4);
    int interactive = 0;
    int debug = 0;

    if (argc < 2) {
        usage(stderr, 1);
```

```
}

int next_option;
do {
    next_option = getopt_long(argc, argv, short_options,
                              long_options, NULL);
    switch(next_option) {
        case 'h':
            usage(stdout, 0);
            // calling usage quits the program anyway but whatever
            // we'll put a break here anyway
            break;

        case 'i':
            interactive = 1;
            break;

        case 'd':
            debug = 1;
            break;

        case '?': // invalid option
            usage(stderr, 1);
            break;

        case -1: // no more options
            break;

        default: // something went very wrong if we got here
            abort();
    }
}
```

```
while (next_option != -1);

// now that we've gone through all the options, OPIND points to first
// nonoption argument. which is hopefully the first filename. unless
// interactive was specified, in which case yell at the user.

if (!interactive) {
    returnargs[2] = string(argv[optind]);
    returnargs[3] = string(argv[optind+1]);
}

if (debug) {
    fprintf(stderr, "Interactive set to %d. Debug is %d.\n", interactive, debug);
    if (!interactive) {
        fprintf(stderr, "Using file %s for graph and %s for subgraph.\n", argv[optind], argv[optind+1]);
    }
}

// put what we want to return in their places.
returnargs[0] = std::to_string(interactive);
returnargs[1] = std::to_string(debug);

return returnargs;
}

void usage (FILE* stream, int exit_code) {
    fprintf(stream, "usage: %s [FILE] [OPTION]\n", program_name);
    fprintf(stream,
        "Options and arguments:\n"
        "FILE:                Space-separated file containing graph adjacency list. If edge weights are specified, they must be in the format 'weight edge'."
        "-h, --help:          Print this help message and exit.\n");
}
```

```
    "-i, --interactive:  Get space-separated graph adjacency list from stdin.\n"
    "-d, --debug:        Print debugging information.\n");

    exit(exit_code);
}

Graph import_data(const char *filename) {

    FILE *fd;

    if (strcmp(filename, "stdin") == 0) {
        fd = stdin;
    }

    else {
        fd = fopen(filename, "r");
        if (fd == NULL) {
            fprintf(stderr, "Could not open file %s for reading.\n", filename);
            perror("Error calling fopen()");
            exit(errno);
        }
    }

    char *line = NULL;
    size_t n = 0;

    Graph g;

    char *tempa, *tempb, *tempweight;
    char tempweight_default[] = "0";
    ssize_t bytes_read;

    // getline returns -1 on failure to read a line (including EOF)
    bytes_read = getline(&line, &n, fd);
```



```
while (bytes_read != -1) {
    tempa = strtok(line, " ");
    // strtok returns NULL if the token (space in this case) was not found in the str
    // otherwise return the string up to the token (exclusive)
    if (tempa == NULL) {
        fprintf(stderr, "Error parsing file %s. Adjacency list must be of form 'node1 node2
            but line was %s\n", filename, line);
        exit(-1);
    }

    // if arg is null after a call to strtok, it'll keep reading after the position of the
    tempb = strtok(NULL, " ");
    if (tempb == NULL) {
        fprintf(stderr, "Error parsing file %s. Adjacency list must be of form 'node1 node2
            but line was %s\n", filename, line);
        exit(-1);
    }

    tempweight = strtok(NULL, " ");
    if (tempweight == NULL) {
        tempweight = tempweight_default;
    }

    int atempa = atoi(tempa);
    int atempb = atoi(tempb);
    int atempweight = atoi(tempweight);

    //printf("Edge is %d %d %d\n", atempa, atempb, atempweight);
    // this will also add the vertices if they don't exist
    // add_vertex is called from within add_edge if no vertex with the given value exists
    g.add_edge(atempa, atempb, atempweight);
    bytes_read = getline(&line, &n, fd);
}
```

```
}

if (fd != stdin) {
    int stat = fclose(fd);
    if (stat != 0) {
        perror("Error closing file");
        // we never wrote to the file so we probably don't even need to check for errors in t
        // but just kill the program if there's an error closing it
        exit(errno);
    }
}

// don't free line since it's null now

return g;
}
```

Listing 8: Code used to check performance of C++ code

References

- Bonnici, V., G., R., Pulvirenti, A., Shasha, D., & Ferro, A. (2013, April). An algorithm for subgraph isomorphism. *BMC Bioinformatics*, 14(Suppl 7)(13). doi:10.1186/1471-2105-14-S7-S13
- Cato, V. (2012, November 24). How to partially compare two graphs. Retrieved November 26, 2017, from <https://stackoverflow.com/a/13537776>
- Cordella, L. P., Foggia, P., Sansone, C., & Vento, M. (2004, October). A (sub)graph isomorphism algorithm for matching large graphs. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 26(10). doi:10.1109/TPAMI.2004.75
- McAuley, J. & Leskovec, J. (2012). Learning to discover social circles in ego networks. *NIPS*. Retrieved November 26, 2017, from <https://snap.stanford.edu/data/egonets-Facebook.html>
- TimeComplexity. (n.d.). Retrieved November 26, 2017, from <https://wiki.python.org/moin/TimeComplexity>
- Ullmann, J. R. (1976, January). An algorithm for subgraph isomorphism. *Journal of the ACM*, 23(1), 31–42. doi:10.1145/321921.321925