

An in-depth study of the syntactic analysis of graphs using  
search matching algorithms, to the degree of sub graphs

Student number: 10637291

email: lectonlm@gmail.com

Supervisor: Dr Linda Marshall

Dept Computer Science, University of Pretoria

August 10, 2015

# 1 Graph Basics

## 1.1 Graph Overview

A Graph in Mathematics and Computer Science is defined as a pair  $G = (V, E)$ , where  $V$  is the set of vertices and  $E$  is the set of edges, formed by pairs of vertices with each other. Figure 1 demonstrates the structural attributes of a simple graph.

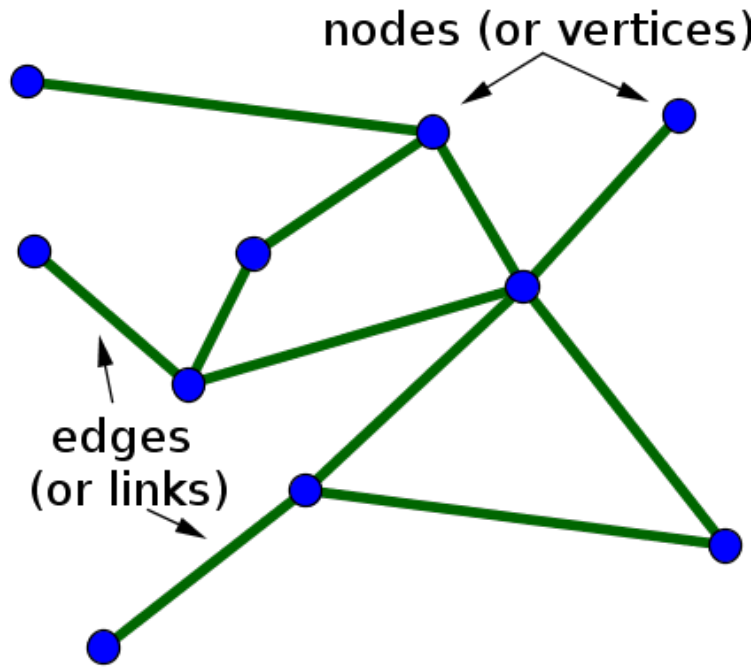


Figure 1: Representation of a graph

Graphs can either be directed or they can be undirected. This means that the edges in the graph could have an absence of direction as in the 1 above, or they could have a direction showing from which vertex the edge is coming from, and to which vertex the edge is going to. The figure below demonstrates a directed graph, commonly known as a digraph. This characteristic is demonstrated by the edge 1, that goes from node a to node b, and also by edges 3 that goes from node c to node a.

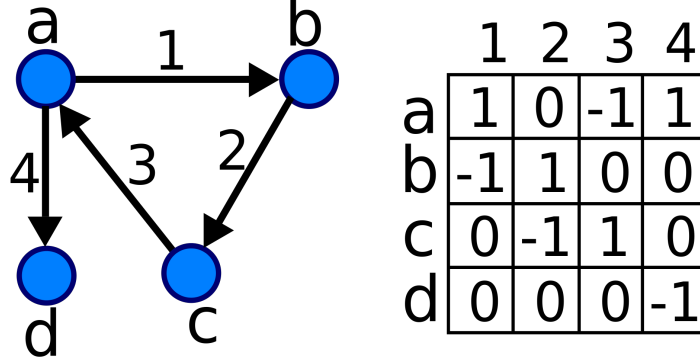


Figure 2: Representation of a digraph

Note that the undirected graph mentioned above is commonly referred to as a bigraph because the direction of its edges could be perceived as to going in both direction as it is not specified.

In this paper, the focus is primarily on directed graphs, and they shall be referred to as digraphs from here onward.

## 1.2 Digraph

A Digraph in mathematics is defined as is a pair of vertices and edges  $(V, E)$  that are disjoint, and their repective mappings that comprise of two components, namely the initial vertex and terminal vertex of each edge i.e. each edge has a initial vertex:

$$E_i \rightarrow V_i \quad (1)$$

and a terminal vertex:

$$V_j \rightarrow E_j \quad (2)$$

for some vertices  $V_i, V_j$  in  $V$  and edges  $E_i, E_j$  in  $E$  [7] refer to the figure 2.

## 1.3 Graph representation

Graphs are represented in a variety of ways, from adjacency lists, incident matrices and adjacency matrices. The algorithms that are studied in this paper make us of adjacency matrices and adjacency list representations of graphs.

### 1.3.1 Adjacency matrices

An adjacency matrices is a  $n \times n$  matrix  $A$ , with  $A(i, j) = 1$  iff  $(i, j) \in E$  [9]. This means that wherever there is an edge in the graph, it is denoted by a 1 in the matrix, places in the matrix  $A$  where there is an absence of an edge  $E$  are denoted by 0.

Figure 3 depicts the association between a graph and its adjacency matrix.

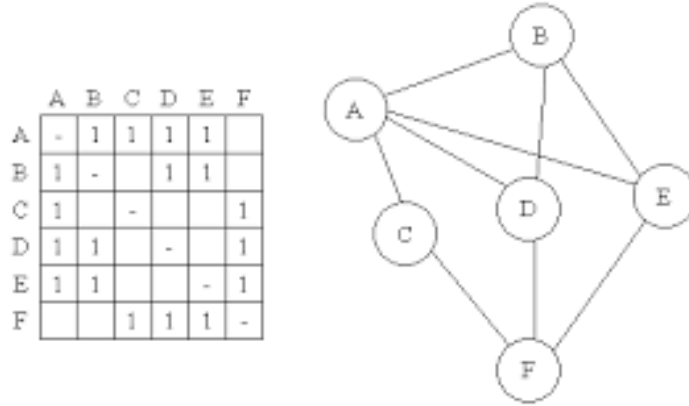


Figure 3: Representation of a graph and its associated adjacency matrix

### 1.3.2 Adjacency list

An Adjacency list is vertices of a graph, of which each vertex is connected to the list. The vertices in an adjacency list point to their own list of edges that they are connected to(i.e. the list contains the edges that connectect them to other vertices). Figure 4 depicts an example of an adjacency list.

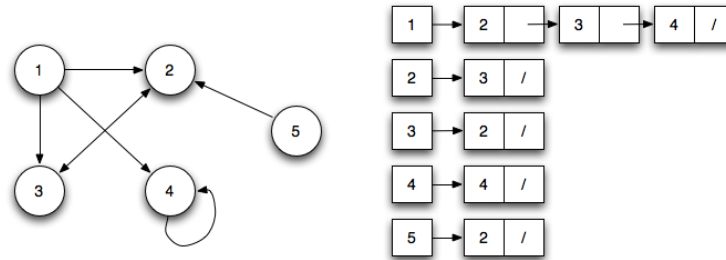


Figure 4: Representation of a graph and its associated adjacency list

## 1.4 Supergraphs and subgraphs

Let  $G_A$  be a graph defined as follows  $G_A = (V_A, E_A)$  and let  $G_B$  be another graph that is defined as follows  $G_B = (V_B, E_B)$  where  $V_A, V_B$  are sets of vertices and  $E_A, E_B$  are sets of edges. In graph theory, a graph  $G_A$  is said to be a subgraph of graph  $G_B$ , and graph  $G_B$  is said to be a supergraph of graph  $G_A$  if all the vertices and edges that are in graph  $G_A$  are also in graph  $G_B$ , that is [3] 1)  $V_A \subseteq V_B$ , and 2) Every edge of  $G_A$  is also an edge in  $G_B$ .

The figure 5 below depicts this relation.

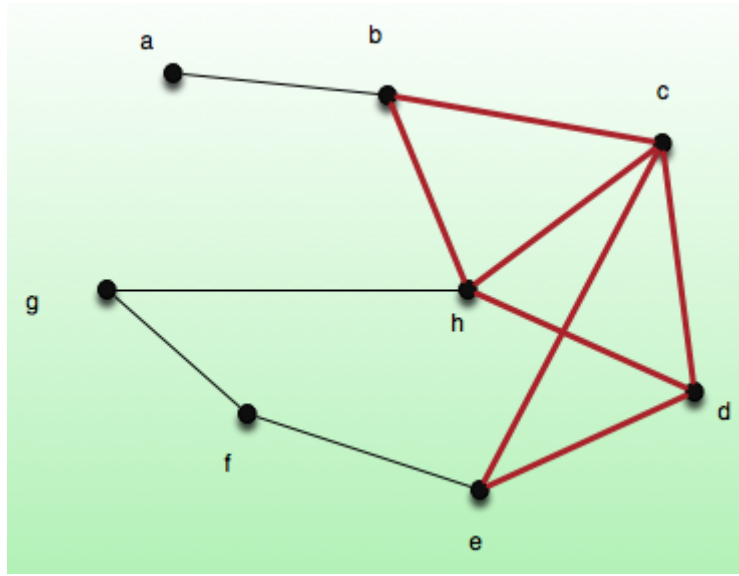


Figure 5: Representation of a super graph and its subgraph that is depicted in red

## 1.5 Graph Isomorphism

Two graphs are said to be isomorphic if they are syntactically similar to each other iff there is a bijection between their respective nodes which make each edge of  $G_1$  correspond to exactly one edge of  $G_2$ , and vice versa[12], i.e. the graphs are structurally the same to each other. This property is demonstrated in figure 6. The two graphs look very different, but when they are further inspected, it is evident that the two are a representation of the same data scheme or even the same graph that has been rearranged.

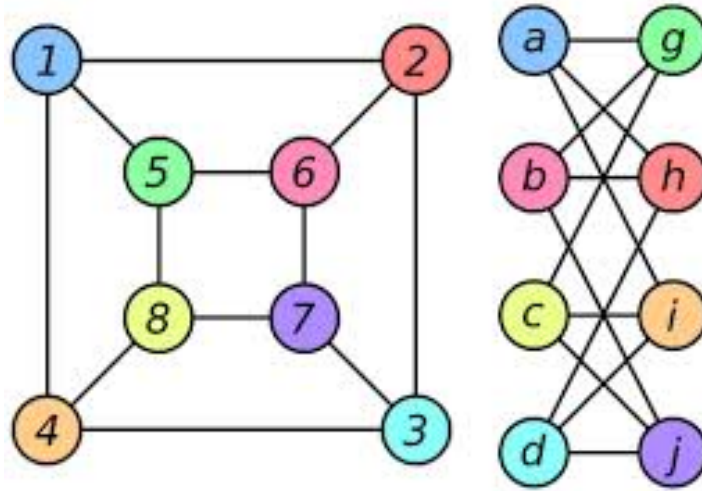


Figure 6: Representation of a super graph and its subgraph that is depicted in red

## 2 Ullman Algorithm

### 2.1 Background

Subgraph isomorphism can be determined using a brute-force approach on the tree representation of a graph  $G_A$ . Though this technique is effective, it is however not efficient because all the possible permutation subgraphs of a graph  $G_A$  are tested against the graph  $G_B$  to determine if there are subgraphs in graph  $G_A$  that are isomorphic to graph  $G_B$ . The number of subgraphs of a graph  $G_A$  increase at an exponential rate with every addition of a vertex  $V_n$  into the graph, thus the total number of subgraphs that are evaluated are

$$ST = 2^{n/2} \quad (3)$$

, where  $ST$  is the total number of subgraphs and  $n$  the number of vertices in the graph  $G_A$ . The matching process is computationally expensive due to this very fact, that is the more vertices and there are in the graph  $G_A$  the more expensive it becomes to detect the subgraph isomorphisms because of the amount of subgraphs it has.

### 2.2 Algorithm

The Ullman algorithm was developed by J.R. Ullman and was published in his titled "An Algorithm for Subgraph Isomorphism" [1]. The algorithm performs graph matching on an adjacency matrix representation of both the graphs, and uses the depth search first (DSF) recursive approach to traverse through the graphs and perform the graph matching process. The Ullman algorithm improves the efficiency of the brute-force approach at detecting subgraph isomorphisms by deductively eliminating nodes in the tree that are in graph  $G_A$ , but are not in graph  $G_B$ , thus reducing the number of subgraphs that are matched against graph  $G_B$  to determine isomorphism.

The algorithm starts by building a starting adjacency matrix  $M_0$  using the two adjacency matrix representations of graphs  $G_A$  and  $G_B$  using the following procedure.

- (1) Construct a  $n \times m$  matrix where  $n$  is the number of rows of graph  $G_B$  and  $m$  is the number of columns of graph  $G_A$ .
- (2) Set all the entries in the matrix to the value of 1.
- (3) Apply the following rule: Set the values in  $M_0$  to 0 for all  $M_{0ij}$  where the degree of a vertex in graph  $G_A$  at  $j$  is greater than the degree of the same vertex in graph  $G_B$  .i.e.

$$deg(A_i) < deg(B_j). \quad (4)$$

A more formal representation of this rule is as follows

$$f(x) = \begin{cases} 1, & \text{if } deg(A_i) \geq deg(B_j) \\ 0, & \text{otherwise } \forall i, j \end{cases}$$

When the starting matrix has been constructed, we systematically permute matrices  $M_d$  from the starting matrix  $M_0$  where  $d$  represents the depth of the generated matrix. The procedure of generating the permuted matrices follows a depth search first (DSF) recursive approach where the stopping condition (leaf matrices) conform to the following form:

- (1) M contains only 0's and 1's.
- (2) There is exactly one 1 in each row.
- (3) Not more than one 1 in each column.

An demonstration of how the permutation matrices are generated is demonstrated in figure 7.

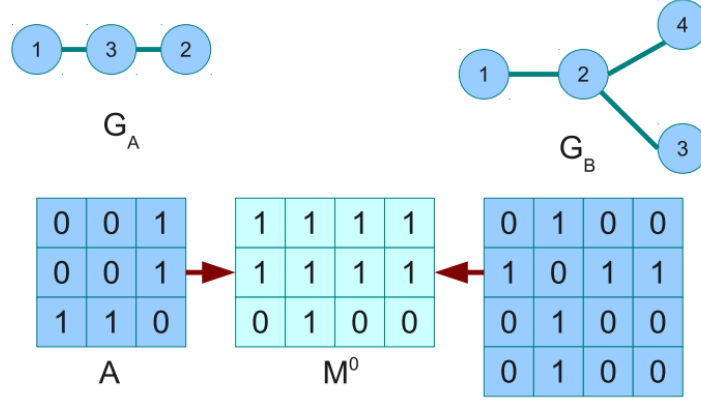


Figure 7: Demonstion of how a permutation matrix is generated from two graphs

Once all the permutation matrices have been generated, each one of the matrices is matched with a graph C, that is obtained from the dot product of the permuted matrix and the graph  $G_A$ . The formula for calculating graph C is follows:  $C = Mn(Mn \cdot G_A)^T$ , where  $G_A$  = input graph A  $Mn$  = permuted matrix Mn in Md, obtained from the starting matrix  $M_0$   $(Mn \cdot G_A)^T$  = the transpose of the dot product of the permuted matrix Mn and the graph  $G_A$ . If there is a single instance of the matrix C, that is calculated using some permuted matrix Mn obtained from the starting matrix  $M_0$ , that is equal to matrix  $G_B$ , then  $G_B$  is isomorphic to  $G_A$ . Thus  $G_B$  is isomorphic to  $G_B$  iff  $G_B$

$$ij = 1 \rightarrow Cij = 1 \forall i, j \quad (5)$$

If non of the generated permuted matrices can statisfy this condition, then  $G_B$  is not isomorphic to  $G_B$ .



## 3 VF2 Algorithm

### 3.1 Background

The VF2 algorithm was introduced by L.P.Cordella, P.Foggiaa, C.Sansone and M.Vento[11]. The algorithm is suitable for graph matching and isomorphic determination, including subgraph isomorphic determination on large graphs, this is attributed to the Data structures that the algorithm uses and the manner in which they are used[11], this feature is discussed later in the paper.

### 3.2 Algorithm

#### 3.2.1 Overview

The algorithm performs the matching process by attempting to find a mapping  $M$ , of vertices in graph  $G_A$  which correspond to vertices in graph  $G_B$ . The mapping is then used to determine if the two graphs are completely syntactically similar(isomorphic), partially syntactically similar or have no structural similarities at all.

#### 3.2.2 Matching and Mapping definition

A mapping  $M$  is defined as a set of pairs  $(n,m)$ , where  $n$  is a vertice from  $G_A$  and  $m$  a vertice from  $G_B$ , thus  $n_j = G_A$  and  $m_j = G_B$ . The isomorphism determining properties of the mapping are defined as follows, a mapping  $M : NA \times NB$  is isomorphic iff  $M$  is a bijection, that preserves the branching structure of  $G_A$  and  $G_B$ , where  $NA$  is a set of vertices from  $G_A$  and  $NB$  a set of vertices from  $G_B$ .

The mapping  $M : NA \times NB$  is subgraph isomorphic iff  $M$  is isomorphic to  $G_A$  and a subgraph of  $G_B$ .

#### 3.2.3 Mapping Procedure

The mapping  $M$  comprises of state based partial solution morphisms  $M(s)$  for each state  $s$ . The process of finding the mapping  $M$  that is described above uses State Space Representations(SSR)[12]. The partial solution morphisms  $M(s)$  selects two subgraphs from  $G_A$  and  $G_B$ , namely  $G_A(s)$  and  $G_B(s)$  respectively. The subgraphs comprises of only vertices that are present in the partial solution  $M(s)$  for the state  $s$  as well as the edges joining them together.

The algorithm starts with an initial state  $s_0$  that has no mapping between the two graphs, thus  $M(s_0) = \text{NULL}$ . The algorithm then computes a set of candidate pairs  $P(s)$ . Each candidate  $p$  in the set is checked against the feasibility function that is discussed in the following chapter, if  $p$  is successful then it is added to the state  $s$ . And the successor  $s'$  is computed using a combination of the predecessor state and the candidate  $p$ , thus:

$$s' = s \cup p \tag{6}$$

The process of generating successor states is a recursive procedure that makes use of the depth first traversal for graphs. When a path has been exhausted and a solution has not yet been found, the algorithm uses backtracking to explore the alternative paths[11,13].

### 3.2.4 Definition of the set $P(s)$ and of the feasibility function $F(s, n, m)$

The VF2 algorithm generates the states with close consideration that only some of the states are consistant with the desired morphisms[12]. The algorithm avoids inconsistant states by making use of a set of rules in it's state generation procedure, thus ensuring that only consistant state are generated, these rules are refered to as feasibility rules.

The algorithm uses a function called a feasibility function to test that an additon of a pair( $n, m$ ) to a state will be consistant. If the addition of the pair passes all the feasibility rules, the algorithm will return a true value, if not, a false value indicating that the procedure results in an inconsistant successor state  $s'$ , and thus that state  $s'$  will not be explored by the algorithm.

A further filter can be applied in the consistent states to rule out those states whose successor states will be inconsistant, this apporoach is employed by adding a additional rules called k-look-ahead rules[12]. They check to see if the current state  $s$  will have a consistant successor state after  $k$  steps, i.e. they check to see if the states from  $s$  to  $s(\text{pow}(k))$  are consistant with the desired morphisms.

### 3.2.5 Condidate Pairs

The candidate pairs are obtained by considering the vertexs that are connect to  $G_A(s)$  and  $G_B(s)$ , the sub-graphs of  $G_A$  and  $G_B$  in the state  $s$ . The vertexs are used to form the pairs  $(n, m)$  as defined above. In order explain how the pairs are formed, we must first introduce the following definitions: Let  $T_{Ain}(s)$  be the set of vertexs that are not yet in the partial mapping  $M(s)$  and are the origin of the edges from graph  $G_A$   $T_{Bin}(s)$  be the set of vertexs that are not yet in the partial mapping  $M(s)$  and are the origin of the edges from graph  $G_B$   $T_{Aout}(s)$  be the set of vertexs that are not yet in the partial mapping  $M(s)$  and are the destination of the edges from graph  $G_A$   $T_{Bout}(s)$  be the set of vertexs that are not yet in the partial mapping  $M(s)$  and are the destination of the edges from graph  $G_B$

The pair  $(n, m)$  is made by vertex  $n$  from  $T_{Aout}(s)$  and  $m$  from  $T_{Bout}(s)$ . If the any of the sets is empty, then we consider the vertex  $n$  from  $T_{Ain}(s)$  and  $m$  from  $T_{Bin}(s)$ . In the case where that graphs are not connected, the pairs will be made by all the vertex not yet contained in either  $G_A(s)$  and  $G_B(s)$ . These pairs form the entries in the set  $P(s)$  for that respective state  $s$ .

### 3.2.6 The feasibility rules

The feasibility rules that are used to ensure that the states that are evaluated play a role in improving the performance, by preventing inconsistent states from being explored and thus optimizing the execution of the algorithm. There are five general feasibility rules defined as  $R_{pred}, R_{succ}, R_{in}, R_{out}$  and  $R_{new}$  respectively A more formal definition is as follows:

$$Fsyn(s, n, m) = Rpred \wedge Rsucc \wedge Rin \wedge Rout \wedge Rnew \quad (7)$$

A more formal definition of each rule is provided below.

The feasibility functions check for two main things, firstly they check the consistency of the partial solution in the successor state  $s'$ , namely  $M(s')$ . Rules  $R_{pred}$  and  $R_{succ}$  are the rules used for this checking. The remaining rules are used for pruning the search space for different levels of look

ahead. The  $R_{in}$  and  $R_{out}$  are used to look ahead one level and determine which of those successor states are consistent, and  $R_{new}$  is used for the same purpose, but for a look ahead level of two. This brings us to the conclusion that, in order for a state to be considered consistent it must pass a combination of all of the five rules, namely:

$$\begin{aligned}
R_{pred}(s, n, m) \iff & \\
& (\forall n' \in M1(s) \cap Pred(GB, n) \exists m' \in Pred(GA, m) | (n', m') \exists M(s) \wedge \\
& \forall m' \in M2(s) \cap Pred(GA, m) \exists n' \in Pred(GB, n) | (n', m') \exists M(s)), \quad (8)
\end{aligned}$$

$$\begin{aligned}
R_{succ}(s, n, m) \iff & \\
& (\forall n' \in M1(s) \cap Succ(GB, n) \exists m' \in Succ(GA, m) | (n', m') \exists M(s) \wedge \\
& \forall m' \in M2(s) \cap Succ(GA, m) \exists n' \in Succ(GB, n) | (n', m') \exists M(s)), \quad (9)
\end{aligned}$$

$$\begin{aligned}
R_{in}(s, n, m) \iff & \\
& (Card(Succ(GB, n) \cap TBin(s)) \geq Card(Succ(GA, m) \cap TAin(s))) \cap \\
& (Card(Pred(GB, n) \cap TBin(s)) \geq Card(Pred(GA, m) \cap TAin(s))) \quad (10)
\end{aligned}$$

$$\begin{aligned}
R_{out}(s, n, m) \iff & \\
& (Card(Succ(GB, n) \cap TBout(s)) \geq Card(Succ(GA, m) \cap TAout(s))) \cap \\
& (Card(Pred(GB, n) \cap TBout(s)) \geq Card(Pred(GA, m) \cap TAout(s))) \quad (11)
\end{aligned}$$

$$\begin{aligned}
R_{new}(s, n, m) \iff & \\
& Card(\tilde{N}A(s) \cap Pred(GB, n)) \geq Card(\tilde{N}B(s) \cap Pred(GA, n)) \wedge \\
& Card(\tilde{N}A(s) \cap Succ(GB, n)) \geq Card(\tilde{N}B(s) \cap Succ(GA, n)) \quad (12)
\end{aligned}$$

This, as described earlier is defined as the feasibility function.