

Simple Necessary Conditions for the Existence of a Hamiltonian Path with Applications to Cactus Graphs

Pascal Welke

Informatik III, University of Bonn, Germany

Abstract

We describe some necessary conditions for the existence of a Hamiltonian path in any graph (in other words, for a graph to be traceable). These conditions result in a linear time algorithm to decide the Hamiltonian path problem for cactus graphs. We apply this algorithm to several molecular databases to report the numbers of graphs that are traceable cactus graphs.

1 Introduction

A Hamiltonian path is a path in a graph G that contains each vertex of G exactly once. The Hamiltonian path problem (i.e., does there exist a Hamiltonian path in a given graph G ?) is a well studied **NP**-complete problem with various applications [5]. Several algorithms have been proposed to find a Hamiltonian path in a graph, or to decide that none exists. For example, Held and Karp [6] give a $O(n^2 \cdot 2^n)$ algorithm to compute a Hamiltonian path. Björklund [1] gives a $O(1.657^n)$ time algorithm to count the number of Hamiltonian paths in a graph, which can also be used to decide the Hamiltonian path problem. Due to the exponential time complexity of those and other algorithms, it would be beneficial to derive simple, fast tests that can be run in advance to decide at least in some cases if there exists a Hamiltonian path, or not.

Many authors concentrated on sufficient conditions for a graph to be traceable (i.e., that it contains a Hamiltonian path). E.g. Dirac [4] gives a lower bound on the number of edges in a graph that implies the existence of a Hamiltonian path. Also, there is a wide range of graph classes, where we know that a Hamiltonian path exists, e.g. complete graphs, cycles, paths, or graphs of the platonic solids.

We go a different way and consider situations which do not allow for a Hamiltonian path. That is, we define easily verifiable properties of graphs that prove that a graph is not traceable. To our knowledge, there is much less work in this direction, most notably by Chvátal [2], that introduces weakly Hamiltonian graphs and derives necessary conditions for a graph to contain a Hamiltonian cycle. However, the paper uses quite involved concepts and the verification of the conditions for a given graph is not straightforward. Our conditions, on the other hand, can be checked in linear time and are easy to understand. They are based on partitioning a graph G into its biconnected components and deriving a tree structure from those objects. In short, a Hamiltonian path in G can only exist if this tree structure is a path.

We start by considering trees and continue by defining a tree structure using the biconnected components of an arbitrary graph to devise conditions in Lemmas 2.3 and 2.4. As a direct application of our necessary conditions, we devise a linear time algorithm for cactus graphs in Theorem 3.1. Finally, we give statistics of a molecular dataset that were obtained using our conditions.

2 Nice Necessary Conditions

From now on, we only consider connected graphs, as otherwise there cannot be a Hamiltonian path. We start by considering the Hamiltonian Path problem for trees. It is easy to see, that a tree T has a Hamiltonian Path if and only if T is a path. We use standard graph notation, see Diestels book [3] for definitions.

Lemma 2.1. *A tree T has a Hamiltonian path if and only if T is a path.*

Proof. “ \Leftarrow ” is clear. “ \Rightarrow ” Let T be a tree and P a Hamiltonian path in T . P contains all vertices of T and has thus $|V(G)| - 1$ edges. Therefore, $E(T) = E(P)$ and thus T is a path. \square

We will show that a generalized version of this holds for a tree structure defined on the *articulation* vertices of any graph G . We need the following definition:

Definition 2.2. *Let G be a connected graph. A vertex $v \in V(G)$ is called articulation vertex if its removal disconnects G , i.e., the graph $G - v = (V', E')$ is disconnected, where $V' := V \setminus \{v\}$ and $E' := \{e \in E : v \notin e\}$. The criticality of v is the number of connected components of $G - v$.*

In a tree, every vertex that is not a leaf is an articulation vertex. We now prove the first necessary condition. In the case of trees, it follows directly from Lemma 2.1.

Lemma 2.3. *Let G be a traceable graph. Then all vertices have criticality at most 2.*

Proof. Suppose there is a vertex v with criticality at least 3. Then $G - v$ has three nonempty connected components C_1, C_2, C_3 . Let P be a Hamiltonian path of G and u_1 (resp. u_2, u_3) be the first vertex in $V(C_1)$ (resp. $V(C_2), V(C_3)$) occurring in P (w.l.o.g. in this order). Any path connecting $u_1 \in V(C_1)$ to $u_2 \in V(C_2)$ in G needs to contain v . Otherwise, u and w would be contained in the same connected component of $G - v$. The same is true for a path from u_2 to u_3 . Therefore, P contains v at least twice, which is a contradiction to P being a path. \square

Figure 1 shows an illustration of the situation described in Lemma 2.3. Vertex v_2 has criticality 3 and therefore does not allow for a Hamiltonian path in the graph. The next lemma focuses on biconnected components.

Lemma 2.4. *Let G be a traceable graph. Then each biconnected component of G contains at most two articulation vertices.*

Proof. Suppose there is a biconnected component B of G that contains three articulation vertices v_1, v_2, v_3 . Removing $v_i \in \{v_1, v_2, v_3\}$ from G results in a disconnected graph $G_i := G - v_i$. Now, there exists a connected component B_i in G_i such that $V(B_i) \cap V(B) = V(B) \setminus \{v_i\}$ and B_i is connected. Let X_i be the nonempty graph of all other connected components of $G - v_i$. To see this, remember that B is a biconnected component thus removing a single vertex does not disconnect B and all vertices in $V(B) \setminus \{v_i\}$ are contained in the same connected component of $G - v_i$. However, as v_i is an articulation vertex, $G - v_i$ is disconnected and thus $V(X_i) \neq \emptyset$. As an example, Figure 1 shows B_i and X_i for the case $v_i = v_1$.

Claim: $V(X_i) \cap V(X_j) = \emptyset$ for all $i \neq j \in \{1, 2, 3\}$.

Using this claim, we can prove the lemma. A Hamiltonian path P of G needs to contain all vertices in $V(X_1), V(X_2), V(X_3)$. But to get from any vertex in $V(X_i)$ to a vertex $x \in V(X_j)$, it needs to pass through v_i . To get from v_i to x , the path must pass through v_j , as $v_i \in V(B_j)$. Using the same argument as in the proof of Lemma 2.3, we see that P needs to visit one of the articulation vertices v_1, v_2, v_3 at least twice, which is a contradiction to P being a path.

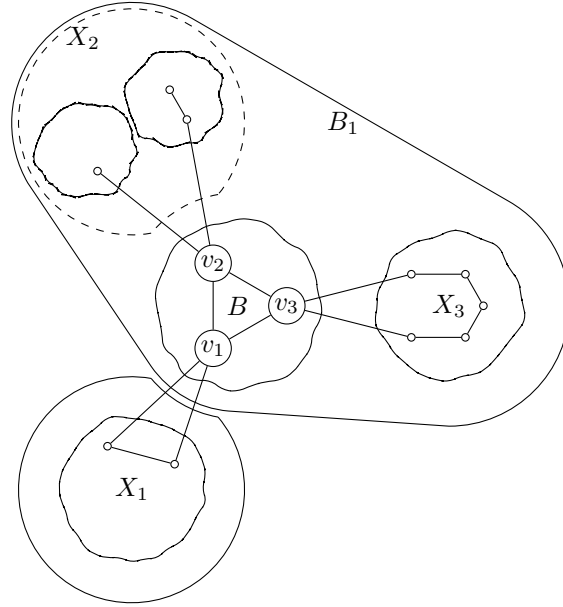


Figure 1: A cactus graph G without a Hamiltonian path. v_2 has criticality 3 (Lemma 2.3) and the biconnected component B contains three articulation vertices (Lemma 2.4).

Proof of Claim: Suppose there exists $x \in V(X_i) \cap V(X_j)$. As $x \in V(X_i)$ there exists a path in X_i connecting x to a neighbor of v_i in G . Thus removing x_j would not disconnect x from $v_i \in V(B_j)$, which contradicts $x \in V(X_j)$. \square

Lemma 2.3 and Lemma 2.4 together show that on any graph G , the existence of a Hamiltonian path implies a path-structure on the biconnected components of G . More exactly, let $\mathcal{A}(G)$ be the set of articulation vertices of G and \mathcal{B} be the set of biconnected components of G . We define a new graph $A(G) = (\mathcal{A}(G), E')$ where $E' = \{\{v, w\} : \exists B \in \mathcal{B} : v, w \in V(B)\}$. Then $A(G)$ is a path. Therefore, the Hamiltonian path problem reduces to checking if these two conditions hold and if there is a Hamiltonian path in each biconnected component, that

- starts at the first articulation vertex and ends at the second articulation vertex (if there are two)
- starts at the articulation vertex (if there is one)
- starts and ends at arbitrary vertices (if there is no articulation vertex in G).

We call biconnected components that contain exactly one articulation vertex *leaf components* and finish this section with an easy corollary of the above considerations.

Corollary 2.5. *Let G be a traceable graph. Then there are either zero or two leaf components.*



Figure 2: A cactus graph on the left and a graph that is not a cactus on the right.

3 The Hamiltonian Path Algorithm for Cactus Graphs

The results of Section 2 imply a polynomial time algorithm for the Hamiltonian path problem for cactus graphs. A *cactus graph* is a connected graph where every biconnected component is either a single edge or a simple cycle. Figure 2 shows a cactus graph and a graph that is no cactus.

Theorem 3.1. *A cactus graph is traceable if and only if all of the following three conditions hold:*

- *Each vertex has criticality at most two*
- *Each biconnected component contains at most two articulation vertices*
- *If a biconnected component contains two articulation vertices, they share an edge.*

Proof. Each cycle is traceable, and each Hamiltonian path of a cycle C starts at an arbitrary vertex of C and ends at one of its two neighbors. Edges are also traceable. “ \Rightarrow ” If a cactus graph G is traceable then, by Lemmas 2.3 and 2.4 the first two conditions hold. Let B be a biconnected component of G that contains two articulation vertices. If B is an edge, then the third condition holds trivially. If B is a cycle, then any Hamiltonian path must enter B through one articulation vertex v , leave it through the other w and can never enter B again. Therefore, the path from v to w must be a Hamiltonian path of B and therefore contains all edges in $E(B)$ except one, which must be $\{v, w\}$. “ \Leftarrow ” We construct a Hamiltonian path as follows: If G is biconnected (i.e., it has no articulation vertices), we construct a Hamiltonian path by removing an arbitrary edge. Otherwise, for each cycle, we remove the edge between the two articulation vertices or one of the edges incident to the unique articulation vertex in the cycle. Note that by this, each articulation vertex has degree two in the resulting graph P . As vertices with criticality 0 have degree one or two in G , every vertex in P has degree less than three. We have removed exactly one edge from each cycle of G , thus P contains no cycles and is still connected. Therefore, P is a path. \square

We can check the conditions of Theorem 3.1 in linear time for a graph G as follows: First, we check if G is connected by a simple breadth first search in linear time. Next, we compute the biconnected components of G in linear time using Tarjans algorithm [9]. Having the biconnected components (given as lists of edges), it is easy to compute the criticality of each vertex in G by counting the number of biconnected components each vertex occurs in as an endpoint of at least one edge. Having the criticality of each vertex, we can compute the number of critical vertices per biconnected components by a single pass over its edge list. To check if G is a cactus graph, we test if each biconnected component is either an edge or a simple cycle, which can also be done by a single pass over all edges in a biconnected component. If there are exactly two, by another pass we can check if the component contains an edge that contains both critical vertices. Therefore, the algorithm can be implemented to run in linear time with a small constant.

4 Some Statistics for Molecular Datasets

We implemented some variants of the proposed algorithm and applied them to three well studied molecular datasets.

NCI-HIV consists of almost $43k$ compounds that are annotated with their activity against the human immunodeficiency virus (HIV) provided by the National Cancer institute [7]. We do not consider the annotations here, but merely use the molecular graph representations. The median number of vertices and edges per graph are 41 and 43, respectively. The database consumes $20.1MB$ in our textual file format.

NCI-2012 is a larger set of molecular graphs from the same source [7]. It consists of more than $250k$ graphs with median number of vertices and edges 36 and 37, respectively. The file size is $100.3MB$.

ZINC is a subset of almost 9 million so called 'Lead-Like' molecules from the zinc database of purchasable chemical compounds [8]. The molecules were selected to have molar mass between $250g/mol$ and $350g/mol$ and have median number of vertices and edges 43 and 44, respectively. The file size is roughly $3.8GB$.

Figure 1 shows the number of graphs N , the number of connected cactus graphs C , the number of traceable cactus graphs T , as well as the number of (arbitrary) graphs X that are definitively not traceable. Furthermore, it reports the time t_i needed by our implementation to compute value $i \in \{N, C, T, X\}$. The numbers were computed by parsing the database from a text file and checking property i for each graph in the respective database. Times were measured using the GNU `time` command summing up `sys` and `user` times.

All experiments were done on an Intel Core i7-2600K with 8GB main memory running Ubuntu 14.04 64bit. The algorithms were implemented in C and compiled using `gcc 4.8.2` with optimization flag `-O3` enabled. No multi-threading was used. Furthermore, due to the fact that each graph can be processed separately, the maximum memory consumption at any time was less than $10MB$.

t_N reports the time our implementation needs to parse the graph database, create graph objects in memory, and dump them again. As you can see, the actual tests only add a small overhead in time compared to just parsing the data. On the other hand, by checking if a graph (a) is connected and (b) fulfills our two necessary conditions, we can declare most of the graphs from all databases as non traceable. For example, for the ZINC dataset, we would only need to further investigate 7 out of almost 9 million graphs to check if they are traceable, or not.

5 Conclusion

We have proposed two necessary conditions for a graph to be traceable that are easy and fast to check. Using them, we proposed a linear time algorithm that decides if a cactus graph is traceable. In more general practical settings, checking these conditions could be a first step, that might, in many cases make applying one of the exponential time exact algorithms obsolete. We evaluated our tests effectiveness in that respect on three molecular data sets of varying size and showed that most molecular graphs can be easily identified as non-traceable using our conditions.

In future work, our proposed algorithm can be extended to yield an exact polynomial time algorithm for more general classes of graphs. Using our conditions, we can reduce the Hamiltonian path problem in a non-biconnected graph G to smaller Hamiltonian path problems in the

	NCI-HIV	NCI-2012	ZINC
N	42687	249533	8946757
C	18028	134478	6517109
X	42658	249436	8946750
T	6	80	0
t_N	0.20	1.01	39.37
t_C	0.28	1.47	56.50
t_X	0.31	1.57	63.94
t_T	0.32	1.67	70.31

Table 1: Statistics for three molecular data sets. The reported times are in seconds

biconnected components of G . We would only need to check if there is a Hamiltonian path in each biconnected component that connects the two articulation vertices or starts at the unique articulation vertex, respectively. This is possible in polynomial time if, for example, the number of spanning trees in each biconnected component is bounded by a polynomial in the size of G .

References

- [1] Andreas Björklund. Determinant sums for undirected hamiltonicity. *SIAM Journal on Computing*, 43(1):280–299, 2014.
- [2] Václav Chvátal. Edmonds polytopes and weakly hamiltonian graphs. *Mathematical Programming*, 5(1):29–40, 1973.
- [3] Reinhard Diestel. *Graph Theory*, volume 173. Springer, 2012.
- [4] GA Dirac. Note on hamilton circuits and hamilton paths. *Mathematische Annalen*, 206(2):139–147, 1973.
- [5] Michael R. Garey and David S. Johnson. *Computers and intractability*, volume 174. Freeman New York, 1979.
- [6] Michael Held and Richard M Karp. A dynamic programming approach to sequencing problems. *Journal of the Society for Industrial & Applied Mathematics*, 10(1):196–210, 1962.
- [7] National Cancer Institute. Cadd group chemoinformatics tools and user services, accessed: 2014-08-31. <http://cactus.nci.nih.gov/>.
- [8] John Irwin. Zinc is not commercial, accessed: 2014-8-20. <http://zinc.docking.org/subsets/lead-like>.
- [9] Robert Tarjan. Depth-first search and linear graph algorithms. *SIAM Journal on Computing*, 1(2):146–160, 1972.