

RESEARCH ARTICLE

# Graphettes: Constant-time determination of graphlet and orbit identity including (possibly disconnected) graphlets up to size 8

Adib Hasan<sup>1</sup>, Po-Chien Chung<sup>2</sup>, Wayne Hayes<sup>2\*</sup>

**1** Ananda Mohan College, Mymensingh, Bangladesh, **2** Dept. of Computer Science, University of California Irvine, Irvine, California, United States of America

\* [whayes@uci.edu](mailto:whayes@uci.edu)



## Abstract

*Graphlets* are small connected induced subgraphs of a larger graph  $G$ . Graphlets are now commonly used to quantify local and global topology of networks in the field. Methods exist to exhaustively enumerate all graphlets (and their orbits) in large networks as efficiently as possible using *orbit counting equations*. However, the number of graphlets in  $G$  is exponential in both the number of nodes and edges in  $G$ . Enumerating them all is already unacceptably expensive on existing large networks, and the problem will only get worse as networks continue to grow in size and density. Here we introduce an efficient method designed to aid statistical *sampling* of graphlets up to size  $k = 8$  from a large network. We define *graphettes* as the generalization of graphlets allowing for disconnected graphlets. Given a particular (undirected) graphette  $g$ , we introduce the idea of the *canonical* graphette  $\mathcal{K}(g)$  as a representative member of the isomorphism group  $Iso(g)$  of  $g$ . We compute the mapping  $\mathcal{K}$ , in the form of a lookup table, from all  $2^{k(k-1)/2}$  undirected graphettes  $g$  of size  $k \leq 8$  to their canonical representatives  $\mathcal{K}(g)$ , as well as the permutation that transforms  $g$  to  $\mathcal{K}(g)$ . We also compute all automorphism orbits for each canonical graphette. Thus, given any  $k \leq 8$  nodes in a graph  $G$ , we can in constant time infer which graphette it is, as well as which orbit each of the  $k$  nodes belongs to. Sampling a large number  $N$  of such  $k$ -sets of nodes provides an approximation of both the distribution of graphlets and orbits across  $G$ , and the orbit degree vector at each node.

## OPEN ACCESS

**Citation:** Hasan A, Chung P-C, Hayes W (2017) Graphettes: Constant-time determination of graphlet and orbit identity including (possibly disconnected) graphlets up to size 8. PLoS ONE 12(8): e0181570. <https://doi.org/10.1371/journal.pone.0181570>

**Editor:** Yongtang Shi, Nankai University, CHINA

**Received:** March 1, 2017

**Accepted:** June 23, 2017

**Published:** August 23, 2017

**Copyright:** © 2017 Hasan et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All the code and the data files (the permutation maps, automorphism orbit lists) are available in [www.github.com/Neehan/Faye](https://www.github.com/Neehan/Faye).

**Funding:** The authors received no specific funding for this work.

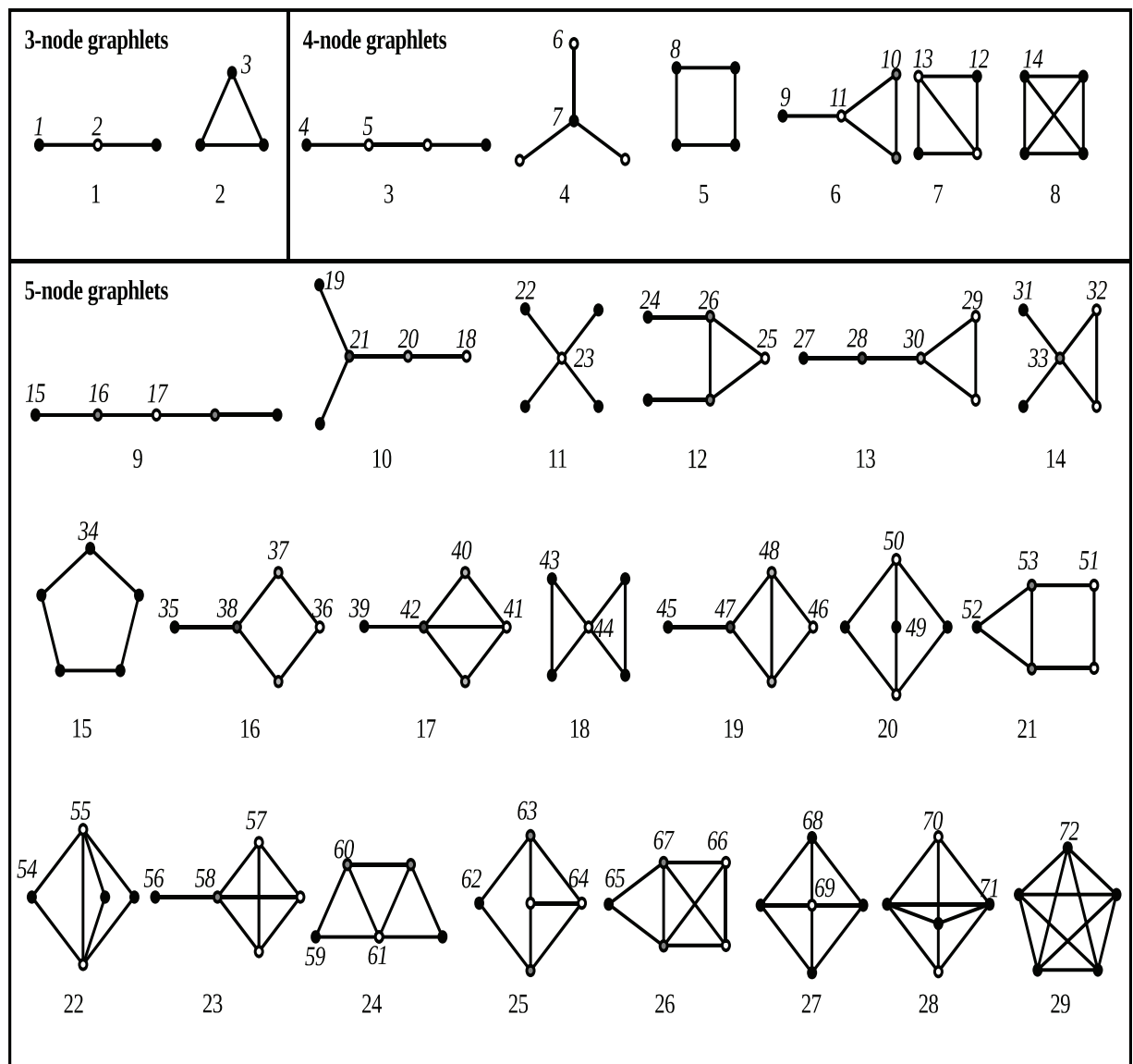
**Competing interests:** The authors have declared that no competing interests exist.

## Introduction

Network comparison is a growing area of research. In general the problem of complete comparison of large networks is intractable, being an  $NP$ -complete problem [1]. Thus, approximate heuristics are needed. Networks have been compared for statistical similarity from a high-level using simple, easy-to-calculate measures such as the degree distribution, clustering coefficients, network centrality, among many others [2, 3]. While more sophisticated methods such as spectral analysis [4, 5] and topological indices [6] have been useful, the study of small subnetworks such as *motifs* [7] and *graphlets* [8, 9] have become popular. They have been used

extensively to globally classify highly disparate types of networks [10] as well as to aid in local measures used to align networks [11–14].

A *graphlet* is a small, connected, induced subgraph  $g$  of a larger graph  $G$ . Given a particular graphlet  $g$ , the *automorphism orbits* of  $g$  are the sets of nodes that are topologically identical to each other inside  $g$ . Graphlets and their automorphism orbits with up to  $k = 5$  nodes were first introduced in 2004 [8], and are depicted in Fig 1. Recently, automated methods have been created that can enumerate, in a larger graph, all graphlets and their automorphism orbits up to graphlet size  $k = 5$  [15] and subsequently to any  $k$  [16], although the latter authors only applied it up to  $k = 6$ . Unfortunately, we have found that these methods take a very long time (hours to



**Fig 1. All (connected) graphlets of sizes  $k = 3, 4, 5$  nodes, and their automorphism orbits; within each graphlet, nodes of equal shading are in the same orbit.** The numbering of these graphlets and orbits were created by hand [8] and do not correspond to the automatically generated numbering used in this paper. The figure is taken verbatim from [16].

<https://doi.org/10.1371/journal.pone.0181570.g001>

days) even just to count graphlets up to size  $k = 5$  on some large biological networks, such as those in BioGRID [17]. It is not clear that such methods, especially for even larger  $k$ , will be applicable to the coming age of ever bigger networks, since the total number of graphlets appearing in a large network tends to increase exponentially with both  $k$  (the graphlet size) and  $n$  (the number of nodes in the large network). Eventually, an exhaustive enumeration of all graphlets appearing in a large network may become infeasible simply due to the number of graphlets that need to be enumerated, even under the optimization of using orbit counting equations. On the other hand, graphlets are too useful to abandon as a method of quantifying the topological structure of graphs. An achievable alternative for a large network  $G$  is to statistically sample its graphlets rather than exhaustively enumerate them. Additionally, such sampling could be useful with the recent advent of comprehensive biological network *databases* [18]: each sampled graphlet would act as a seed for local matching between larger networks, similar to how  $k$ -mers (short sequences of length  $k$ ) are used for seed-and-extend sequence matching in BLAST [19].

To efficiently create a statistical sample of graphlets in a large network  $G$ , one must be able to take an arbitrary set of  $k$  nodes from  $G$ , and efficiently (preferably in constant time) determine both *which* graphlet is represented, as well as the automorphism orbits of each of the  $k$  nodes. Here, we solve this problem both by enumerating all graphlets (and their disconnected counterparts, which we term *graphettes*) and their automorphism orbits up to graphettes of size  $k = 8$ . We present a method that creates a lookup table that can quickly determine the graphette identity of any  $k$  nodes, as well as their automorphism orbits. Since the lookup table required significant time to pre-compute for  $k = 7$  (a few hours on a single core) and  $k = 8$  (hundreds of CPU weeks on a cluster), we provide the actual lookup tables for these values of  $k$  online at <http://github.com/Neehan/Faye>.

## Materials and methods

### Definitions and notations

Given a graph  $G$  on  $n$  nodes, a  $k$ -*graphette* is a (not necessarily connected) induced subgraph  $g$  on any set of  $k$  nodes of  $G$ . There are many ways one could choose the  $k$  nodes, for example (i) choosing  $k$  nodes uniformly at random from  $G$ , or (ii) performing a local search around some node  $u$ . We expect the former to be useful only in dense networks, while the latter is probably more useful in sparse networks because most random sets of  $k$  nodes in a sparse graph will be highly disconnected and thus not very informative. One could also (iii) perform edge-based selection (with local expansion) to ensure dense regions are sampled more frequently than sparse regions [20]; still other methods have been suggested [21].

Given a set of  $k$  nodes, we wish to quickly ascertain which graphette is represented, and which automorphism orbits each of the  $k$  nodes belong to. To do that we need a canonical list of graphettes and their orbits, and a fast way to determine which canonical graphette is represented by any permutation of  $k$  nodes. Here we demonstrate how, if  $k$  is fixed and relatively small ( $k \leq 8$  in our case), this can be accomplished in constant time by pre-computing and storing a lookup table indexed by a bit vector representation of the lower triangular matrix of the (undirected) adjacency matrix of the induced subgraph. Given such an index, the value associated with that index identifies the canonical graphette (a canonical ordering of the nodes for that graphette). We also pre-compute the automorphism orbits of all the canonical graphettes. Thus, by reversing the lookup table we can, in constant time, infer the orbit identity of each of the  $k$  nodes in that  $k$ -graphette. As a corollary, we can also update the (statistically sampled) *graphette orbit degree vector* of each of the  $k$  nodes, similar to the graphlet degree vector [9].

We use the following abbreviations and notations throughout:

$G(V, E)$	The Graph with nodes $V$ and edges $E$
$\mathcal{V}(G)$	The set of nodes of graph $G$
$\mathcal{E}(G, u, v)$	The boolean value denoting connectivity between nodes $u$ and $v$ of graph $G$
$\iff$ , iff	If and only if
$ S $	The number of elements in set $S$ .
$Adj(G)$	The adjacency matrix representation of graph $G$
$Aut(G)$	The set of automorphisms of graph $G$
$\mathcal{K}(g)$	Canonical isomorph of graphette $g$

<https://doi.org/10.1371/journal.pone.0181570.t001>

### Canonization of graphettes

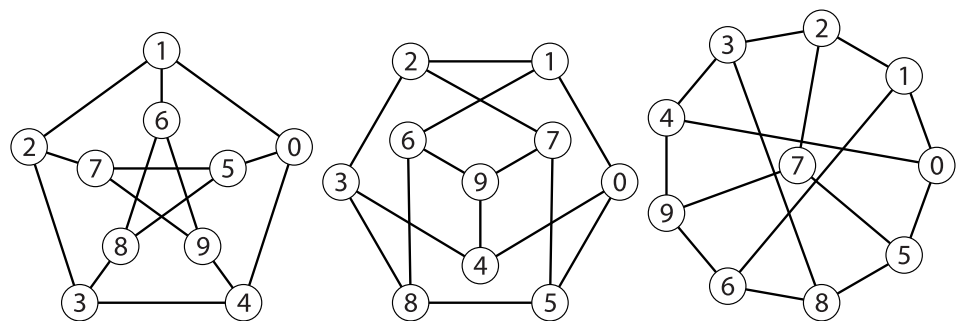
If graphs  $G$  and  $H$  are isomorphic, it essentially means they are exactly the same graph, but drawn differently. For example, Fig 2 shows three different drawings of the Petersen graph. Technically, an isomorphism between networks  $G$  and  $H$  is a permutation  $\pi : \mathcal{V}(G) \rightarrow \mathcal{V}(H)$  so that

$$\mathcal{E}(G, u, v) \iff \mathcal{E}(H, \pi(u), \pi(v)),$$

Consider a 3-graphette with nodes  $w, x$  and  $y$ . There are only 4 possible such graphettes, depicted in Fig 3. However, by permuting the order of the nodes, each of these graphettes can be represented by several isomorphic variants. In order to determine if two graphettes are isomorphic, we will represent its (undirected) graph with the lower-triangle of its adjacency matrix. We will place this lower-triangular matrix into a bit vector, resulting in a representation similar to existing ones for orbit identification [16].

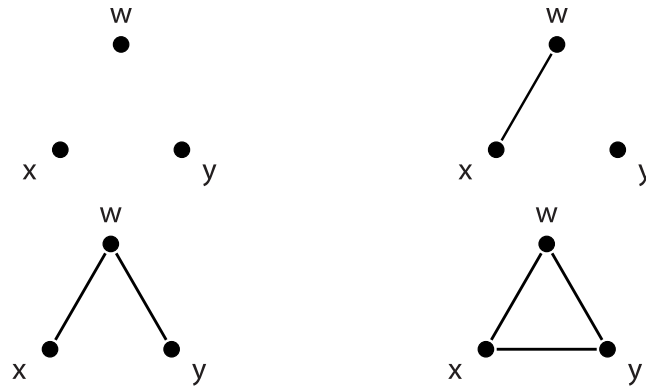
We now describe the idea of a *canonical representative* of each isomorph. To provide an explicit example, consider Fig 4, depicting the three isomorphic configurations of the 3-graphette that has exactly one edge. In order to determine that these graphettes are all isomorphic, we take the bit vector representation depicted, and define the lowest-numbered bit-vector among all the isomorphs as the *canonical representative*. All the other isomorphs in the lookup table point to it. In this way, every graph on 3 nodes can be efficiently mapped to its canonical 3-isomorph.

We also automatically determine the number of *automorphism orbits* (see below) for each canonical isomorph. Table 1 represents, for various values of  $k$ , the number of bits  $b(k)$  required



**Fig 2. Three isomorphic representations of the Petersen graph.**

<https://doi.org/10.1371/journal.pone.0181570.g002>



**Fig 3. All the possible 3-graphettes.**

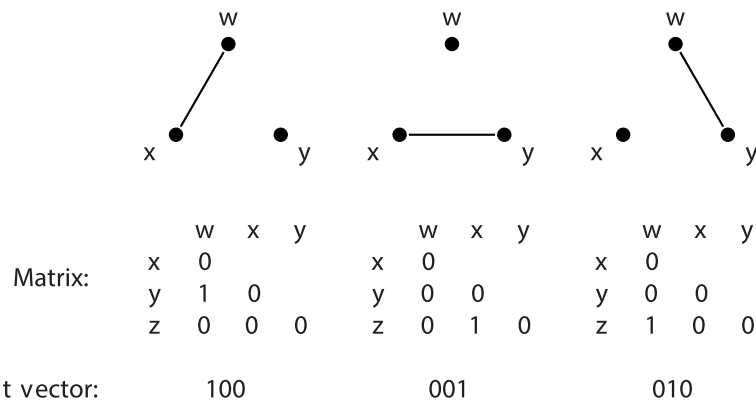
<https://doi.org/10.1371/journal.pone.0181570.g003>

to store the lower-triangular matrix of all graphettes on  $k$  nodes (i.e., the length of the bit vector used to store this matrix); the resulting total number possible representations of  $k$  nodes (which is simply  $2^{b(k)}$ ); the number of canonical isomorphs  $NC(k)$ ; and the number of canonical automorphism orbits. Note that, to map each possible set of  $k$  nodes to their canonical isomorphs, the lookup table has  $2^{b(k)}$  entries, and each entry has a value between 0 and  $NC(k) - 1$ . Note that for  $k$  up to 8, the graphettes can be stored in 32 bits. In that case, the maximum space required will be  $32 \times 2^{28} = 1$  GB. This is as far as we go, for now. Moore’s Law suggests that we may be able to go to  $k = 9$  within a few years, and to  $k = 10$  in perhaps a decade or two.

We note that the most expensive part of our algorithm is creating the lookup table between an arbitrary set of  $k$  nodes, to the canonical graphette represented by those  $k$  nodes; in the absence of a requirement for this lookup table, one could use orbit counting equations [16] to generate automorphism orbits up to  $k = 12$ .

### Generating the lookup table from non-canonical to canonical graphettes

Assume the large graph  $G$  has  $n$  nodes labeled 0 through  $n - 1$ , and pick an arbitrary set of  $k$  nodes  $U = \{u_0, u_1, \dots, u_{k-1}\}$ . Create the subgraph  $g$  induced on the nodes in  $U \subseteq \mathcal{V}(G)$ , and



**Fig 4. All 3-graphettes with exactly one edge; the canonical one is the one with lowest integer representation (the middle one in this case).** Each of them is placed in a lookup table indexed by the bit vector representation of its adjacency matrix, pointing at the canonical one. In this way we can determine that it is the one-edge 3-graphette in constant time.

<https://doi.org/10.1371/journal.pone.0181570.g004>

**Table 1.** For each value of  $k$ : The number of bits  $b(k) = \frac{k(k-1)}{2}$  required to store the lower-triangle of the adjacency matrix for an undirected  $k$ -graphette; the number of such  $k$ -graphettes counting all isomorphs which is just  $2^{b(k)}$ ; the number of canonical  $k$ -graphettes (this will be the number of unique entries in the above lookup table [22], and up to  $k=8$ , 14 bits is sufficient); and the total number of unique automorphism orbits (up to  $k=8$ , 17 bits is sufficient) [27]. Note that up to  $k=8$ , together the lookup table for canonical graphettes and their canonical orbits fits into 31 bits, allowing storage as a single 4-byte integer, with 1 bit to store whether the graphette is connected (i.e., also a graphlet). The suffixes K, M, G, T, P, and E represent exactly  $2^{10}$ ,  $2^{20}$ ,  $2^{30}$ ,  $2^{40}$ ,  $2^{50}$  and  $2^{60}$ , respectively.

$k$	bits $b(k)$	#Graphs $2^{b(k)}$	Space $b(k)2^{b(k)}$	#Canonicals $NC(k)$	#Orbits
1	0	1	0	1	1
2	1	2	0.25 B	2	2
3	3	8	3 B	4	6
4	6	64	48 B	11	20
5	10	1 K	1.25 KB	34	90
6	15	32 K	60 KB	156	544
7	21	2 M	5.25 MB	1044	5096
<b>8</b>	<b>28</b>	<b>256 M</b>	<b>896 MB</b>	<b>12346</b>	<b>79264</b>
9	36	64 G	288 GB	274668	2208612
10	45	32 T	180 TB	12005168	113743760
11	55	32 P	220 PB	1018997864	10926227136
12	66	64 E	528 EB	165091172592	1956363435360

<https://doi.org/10.1371/journal.pone.0181570.t002>

let its bit vector representation  $B$  be of the form lower-triangular matrix described in Fig 4. We now describe how to create the lookup table that maps any such  $B$  to its canonical representative.

We iterate through all  $2^{b(k)}$  bit vectors in order; for each value  $B$ , we check to see if it is isomorphic to any of the previously found canonical graphettes; if so, the lookup table value is set to the previously found canonical graphette; otherwise we have a new, previously unseen canonical graphette and the lookup table value is set to itself ( $B$ ).

When checking for isomorphism between  $B$  and all previously found canonical graphettes, we use a relatively simple brute force approach. If the degree distribution of the two graphettes are different, we can immediately discard the pair as non-isomorphic; otherwise we resort to cycling through every permutation of the nodes checking each pair for graph equality, which has worst-case running time of  $k^2k!$ . The total run time to compute the lookup table for a particular value  $k$  is thus bounded above by  $k^2k! \cdot NC(k) \cdot 2^{b(k)}$ , where  $k!$  is the maximum number of permutations we need to check if a non-canonical matches an existing canonical,  $k^2$  is the worst-case running time to check if 2 specific permutations of  $k$ -graphettes are isomorphic, there are at most  $NC(k)$  canonicals to check against [22], and  $2^{b(k)} = 2^{n(n-1)/2}$  is the total number of undirected graphs on  $k$  nodes. More sophisticated approaches exist [23], which may more easily allow higher values of  $k$ .

This process can also be parallelized, which is what we did for  $k=8$ . Essentially, we can split the  $2^{b(k)}$  non-canonical graphettes into  $m$  sets of about  $2^{b(k)}/m$  graphettes each, and then spread the computation across  $m$  machines. For each of the  $m$  sets  $S_i$ , we loop through all graphettes in that set and mark out which are isomorphic to each other. For each set  $S_i$ , we will find a set  $T_i$  of lowest-numbered “temporary” canonical graphettes in  $S_i$ , along with the map  $TC: S_i \rightarrow T_i$  of which graphettes in  $S_i$  map to each temporary canonical in  $T_i$ . That is, for each graphette  $g \in S_i$ ,  $\exists h \in T_i$  for which the temporary canonical  $TC(g) = h$ . Finally, once all the  $m$  sets have been evaluated in this way, a second stage passes through all the  $T_i$ ,  $i = 0, \dots, m - 1$ , merging the temporary canonicals together into a final, global list of canonical graphettes, while also propagating these globally lowest-numbered canonicals back up through the  $m$  temporary canonical

maps, so each graphette  $g$  globally maps to the globally lowest-numbered canonical; we call this process *sifting for canonicals*, and it may require several iterations to globally find the final list of canonicals. In this way we ran  $k = 8$  in about a week across 600 cores, for a total of 600 CPU-weeks. This process could probably be made more efficient with smarter isomorphism checking [23, 24].

### Graph automorphism and orbits

An isomorphism  $\pi : \mathcal{V}(g) \rightarrow \mathcal{V}(g)$  (from a graph  $g$  to itself) is called an *automorphism*.

While an isomorphism is just a permutation of the nodes, it is called an *automorphism* if it results in exactly the same labeling of the nodes in the same order—in other words exactly the same adjacency matrix. The set of all automorphisms of  $g$  will be called  $Aut(g)$ .

An *automorphism orbit*, or just *orbit*, of  $g$  is a minimally sized collection of nodes from  $\mathcal{V}(g)$  that remain invariant under **every** automorphism of  $g$  [25]. There can be more than one automorphism orbit, and each orbit can have anywhere from 1 to  $k$  member nodes; refer again to Fig 1 for some examples. More formally, a set of nodes  $\omega$  constitute an orbit of  $g$  iff:

1. For any node  $u \in \omega$  and **any** automorphism  $\pi$  of  $g$ ,  $u \in \omega \iff \pi(u) \in \omega$ .
2. if nodes  $u, v \in \omega$ , then there exists an automorphism  $\pi$  of  $g$  and a  $\gamma > 0$  so that  $\pi^\gamma(u) = v$ .

Now, we shall prove a few relevant results that will be useful later for automatically enumerating the orbits.

**Proposition 1.** *For each node  $u \in \mathcal{V}(g)$  and each automorphism  $\pi : \mathcal{V}(g) \rightarrow \mathcal{V}(g)$ , there exists an integer  $\lambda > 0$  such that  $\pi^\lambda(u) = u$ .*

*Proof.* Because  $\pi$  is an automorphism,

$$\begin{aligned} u \in \mathcal{V}(g) &\implies \pi(u) \in \mathcal{V}(g) \\ &\implies \pi^2(u) \in \mathcal{V}(g) \\ &\vdots \\ &\implies \pi^i(u) \in \mathcal{V}(g), \forall i \in \mathbb{N}. \end{aligned}$$

Since  $|\mathcal{V}(g)|$  is finite and  $\pi$  is bijective, the conclusion obviously follows.

We shall call the set of nodes

$$\mathcal{C}_\pi(u) = \{u, \pi(u), \dots, \pi^{\lambda-1}(u)\}$$

the *cycle* of  $u$  under automorphism  $\pi$ , where  $\lambda$  is the smallest positive integer such that  $\pi^\lambda(u) = u$ .

Note that  $\lambda$  is not unique since  $\pi^\lambda(u) = \pi^{2\lambda}(u) = \dots = u$ . Also,  $\pi$ ,  $u$ , and  $\lambda$  are tied together into triples such that knowing any two determines the third.

**Corollary 1.1.**  *$\pi$  maps every node  $\in \mathcal{C}_\pi(u)$  to a node (possibly same)  $\in \mathcal{C}_\pi(u)$ .*

**Corollary 1.2.** *In any automorphism  $\pi$  of  $g$ , every node appears in exactly one cycle.*

In other words, the cycles  $\pi$  creates are disjoint. (However, the cycles from different automorphisms might not be so.) Hence, it makes sense to say *splitting an automorphism into its cycles*. For example consider the permutation  $\pi = (201354)$  of  $(012345)$ . Since  $\pi(0) = 2$ ,  $\pi(2) = 1$ ,  $\pi(1) = 0$ , the nodes  $(012)$  form a cycle. Now start with the next node, 3.  $\pi(3) = 3$ . So,  $(3)$  is another cycle. Finally,  $\pi(4) = 5$ ,  $\pi(5) = 4$ , so,  $(45)$  form another cycle. Hence, the permutation  $(201354)$  is split into three cycles, namely  $(012)$ ,  $(3)$ ,  $(45)$ .

**Proposition 2.** *The orbits are disjoint. (In other words, each node appears in exactly one orbit.)*

*Proof.* Assume the contrary, i.e., a node  $u \in \mathcal{V}(g)$  appears in two different orbits  $\omega_1$  and  $\omega_2$ . According to the second condition, for any other node  $v \in \omega_1$ , there exists an automorphism  $\pi$  of  $g$  and a  $\gamma$  so that  $\pi^\gamma(u) = v$ . However, from the first condition,

$$\begin{aligned} u \in \omega_2 &\implies \pi(u) \in \omega_2 \\ &\implies \pi^2(u) \in \omega_2 \\ &\quad \vdots \\ &\implies \pi^\gamma(u) \in \omega_2 \\ &\implies v \in \omega_2 \end{aligned}$$

Therefore, every node  $v \in \omega_1$  also belongs to  $\omega_2$ . Hence,  $\omega_1 \subseteq \omega_2$ .

Following the same logic,  $\omega_2 \subseteq \omega_1$ , implying  $\omega_1 = \omega_2$ .  $\Rightarrow \Leftarrow$

**Corollary 2.1.** *Each cycle appears in exactly one orbit, which completely contains that cycle.*

*Proof.* If an orbit  $\omega$  partially contains a cycle  $\mathcal{C}_\pi(u)$ , then  $\omega$  is not invariant under automorphism  $\pi$ , as  $\pi$  will map some node in  $\omega$  (and  $\mathcal{C}_\pi(u)$ ) to another node outside  $\omega$  (but still in  $\mathcal{C}_\pi(u)$ ) according to corollary 1.1, contradicting our definition of orbits. Since two orbits are disjoint,  $\mathcal{C}_\pi(u)$  must appear only in  $\omega$ , and in none of the other orbits.

These statements are enough to be able to find all orbits of each graphette, as we now demonstrate.

## Automatically enumerating all orbits of a graph

From the propositions in the previous section, an algorithm to enumerate the orbits can be constructed like this:

1. Generate all automorphisms of  $g$ .
2. Split each automorphism into its cycles.
3. Merge the cycles from different automorphisms to form orbits.

**Generating all automorphisms of  $g$ .** Referring to Algorithm 1, the function `GENERATEAUTOMORPHISMS()` applies every possible permutation of  $\mathcal{V}(g)$  over  $Adj(g)$ . Each permutation creates an isomorph of  $Adj(g)$ . If  $Adj(g)$  is unchanged under some permutation  $\pi$ , then by definition,  $\pi$  is an automorphism of  $g$ . Hence it is saved into  $Aut(g)$ .

Two optimization strategies are employed:

1. No node is mapped to another node with unequal degree.
2. An automorphism of graph  $g$  is also an automorphism of its complement graph  $g'$ .

In practice, this algorithm generates all automorphisms of all the canonical graphettes up to size 8 in a matter of seconds. Nevertheless, for additional speed up in higher sizes, modern sophisticated automorphism detection algorithms [23, 24] may be used.

**Splitting automorphisms into cycles.** An automorphism  $\pi$  of  $g$  is basically a permutation of nodes of  $g$ . Hence, to split  $\pi$  into cycles, we can repeatedly apply  $\pi$  over every node  $u \in \pi$  and remember the nodes  $u$  transforms into. This forms the cycle with node  $u$ , i.e.  $\mathcal{C}_\pi(u)$ , which is saved in  $\mathcal{C}$ . After first visit, each node is marked visited to prevent more visits.

**Merging cycles to enumerate orbits.** Suppose  $\mathcal{C}(g)$  is the set of all cycles resulting from all the automorphisms of  $g$ .



To enumerate orbits from it, first each node  $u$  is colored with a unique color  $\omega(u) = u$ . Then  $\omega(u)$  is continuously updated to reflect the current color of  $u$ , as the nodes belonging to same orbits are gradually colored by identical color.

For the nodes of each cycle  $c \in \mathcal{C}(g)$ , we save their minimum color in  $\omega_{\min}$ , and then color all of them with  $\omega_{\min}$ . After coloring all the cycles in this way, nodes belonging to same orbits get the same color, and hence, get enumerated.

**Algorithm 1** Automatically enumerating automorphism orbits of a graph

```

function GENERATEAUTOMORPHISMS (Graph  $g$ )
   $Aut(g) = \{\}$  // Find the automorphisms of  $g$ 
  for each permutation  $\pi$  of  $\mathcal{V}(g)$  do
    apply  $\pi$  over  $Adj(g)$ 
    if  $Adj(g) == \pi(Adj(g))$  then put  $\pi$  in  $Aut(g)$ 
    end if
  end for
end function
function GENERATECYCLES (automorphism  $\pi$ )
   $\mathcal{C} = \{\}$ 
  for node  $u$  in  $\pi$  do
    if  $u$  is not visited then
      mark  $u$  visited
      new cycle  $\mathcal{C}_\pi(u) = \{\}$ 
      node  $v = \pi(u)$ 
      while  $v \neq u$  do
        put  $v$  in  $\mathcal{C}_\pi(u)$ 
        mark  $v$  visited
         $v = \pi(v)$ 
      end while
      put  $\mathcal{C}_\pi(u)$  in  $\mathcal{C}$ 
    end if
  end for
end function
function ENUMERATEORBITS ( $\mathcal{C}(g)$ )
  for each node  $u \in \mathcal{V}(g)$  do  $\omega(u) = u$ 
  end for
  for cycle  $c \in \mathcal{C}(g)$  do
    let  $\omega_{\min} = \infty$ 
    for node  $u \in c$  do  $\omega_{\min} = \min(\omega_{\min}, \omega(u))$ 
    end for
    for node  $u \in c$  do  $\omega(u) = \omega_{\min}$ 
    end for
  end for
end function

```

**Proof of correctness of Algorithm 1**

Here we prove that Algorithm 1 determines every orbit of  $g$ .

Suppose a set  $\omega$  is among the final sets generated by Algorithm 1. We shall prove  $\omega$  is an orbit of  $g$  by showing that it follows the two properties of orbits:

1. Let a node  $u \in \omega$  form the cycle  $\mathcal{C}_\pi(u)$  under automorphism  $\pi$ . The GENERATECYCLES function will apply  $\pi$  repeatedly until it finds a  $\lambda$  so that  $\pi^\lambda(u) = u$  and will therefore determine  $\mathcal{C}_\pi(u)$ . Since the ENUMERATEORBITS function assigned  $u$  to  $\omega$ , it had also assigned all nodes in  $\mathcal{C}_\pi(u)$  to  $\omega$ . Hence  $u \in \omega \iff \pi(u) \in \omega$ .
2. Suppose nodes  $u, v \in \omega$ . Then, either they belonged to a cycle from which they were assigned to a mutual set  $\omega$  in ENUMERATEORBITS function, or there is a third node  $w$  so that

$w$  shares separate cycles with  $u$  and  $v$  under different automorphisms  $\pi_1$  and  $\pi_2$ . In the first case,  $u$  and  $v$  already belong to a common cycle. In the second case, assume  $\pi_1^{\gamma_1}(w) = u$  and  $\pi_2^{\gamma_2}(w) = v$ . Consider the permutation  $\phi = \pi_2^{\gamma_2} \circ \pi_1^{-\gamma_1}$ . Since composition of two automorphisms is an automorphism [26],  $\phi$  is also an automorphism. And notice that

$$\phi(u) = \pi_2^{\gamma_2}(\pi_1^{-\gamma_1}(u)) = \pi_2^{\gamma_2}(w) = v$$

implying  $u$  and  $v$  belong to a common cycle under  $\phi$ .

Therefore,  $\omega$  is indeed an orbit of  $g$ . Since each node was given a unique orbit color in the beginning of `ENUMERATEORBITS`, every orbit of  $g$  will be eventually found by Algorithm 1.

## Results and discussion

Using the algorithms described herein, we have enumerated all possible graphlets, including the generalization of disconnected counterparts called *graphettes*, up to size  $k = 8$ . The code and data can be found in <http://github.com/Neehan/Faye>. (Note that the github code uses the *upper* triangle matrix, though we intend to convert it to use the lower triangle as that representation has already been established [16].) We have also enumerated all orbits up to size  $k = 8$ . More importantly to the statistical sampling technique described in the Introduction, we have used a bit-vector representation of all possible adjacency matrices of all possible sets of up to  $k = 8$  nodes and created a lookup table from the  $2^{k(k-1)/2}$   $k$ -sets to their canonical graphette representatives. This allows us to determine, in constant time, the graphette represented by these  $k$  nodes, as well as the automorphism orbits of each nodes. This allows efficient estimation of both the global distribution of graphlets and orbits, as well as an estimation of the graphlet (or orbit) degree vector for each node in a large graph  $G$ .

Although the lookup tables for  $k > 8$  are at present too big to compute or store, we could also use NAUTY or SAUCY to enumerate all the canonical graphettes up to size  $k = 12$ , and use our orbit generation code Algorithm 1 to determine all the orbits in all graphettes up to size  $k = 12$ . We have verified that previous results are consistent with ours in terms of the number of distinct graphettes [22] and orbits [27] determined, as displayed in Table 1.

In future work we will study which statistical sampling techniques most efficiently produce a good estimate of the complete graphlet and local (per-node) degree vectors. We also intend to study how this method may aid in cataloging of graphlets for database network queries, or in non-alignment network comparison [10]. Finally, there may be ways to combine our method with those of orbit counting equations [15, 16] to more efficiently produce samples of orbit counts.

## Acknowledgments

We thank Sridevi Maharaj, Dillon Kanne, and the anonymous referees for several helpful suggestions on presentation.

## Author Contributions

**Conceptualization:** Adib Hasan, Wayne Hayes.

**Formal analysis:** Adib Hasan, Wayne Hayes.

**Investigation:** Po-Chien Chung, Wayne Hayes.

**Methodology:** Po-Chien Chung, Wayne Hayes.

**Project administration:** Wayne Hayes.

**Resources:** Wayne Hayes.

**Software:** Adib Hasan, Po-Chien Chung.

**Supervision:** Wayne Hayes.

**Validation:** Adib Hasan.

**Writing – original draft:** Adib Hasan, Wayne Hayes.

**Writing – review & editing:** Adib Hasan, Wayne Hayes.

## References

1. Cook SA. The Complexity of Theorem-proving Procedures. In: Proceedings of the Third Annual ACM Symposium on Theory of Computing. STOC'71. New York, NY, USA: ACM; 1971. p. 151–158. Available from: <http://doi.acm.org/10.1145/800157.805047>.
2. Newman M. Networks: an introduction. 2010. United States: Oxford University Press Inc, New York. 2010; p. 1–2.
3. Emmert-Streib F, Dehmer M, Shi Y. Fifty years of graph matching, network alignment and network comparison. *Information Sciences*. 2016; 346:180–197. <https://doi.org/10.1016/j.ins.2016.01.074>
4. Wilson RC, Zhu P. A study of graph spectra for comparing graphs and trees. *Pattern Recognition*. 2008; 41(9):2833–2841. <https://doi.org/10.1016/j.patcog.2008.03.011>
5. Thorne T, Stumpf MP. Graph spectral analysis of protein interaction network evolution. *Journal of The Royal Society Interface*. 2012; p. rsif20120220.
6. Dehmer M, Emmert-Streib F, Shi Y. Interrelations of graph distance measures based on topological indices. *PLoS one*. 2014; 9(4):e94985. <https://doi.org/10.1371/journal.pone.0094985> PMID: 24759679
7. Milo R, Shen-Orr S, Itzkovitz S, Kashtan N, Chklovskii D, Alon U. Network motifs: simple building blocks of complex networks. *Science*. 2002; 298(5594):824–827. <https://doi.org/10.1126/science.298.5594.824> PMID: 12399590
8. Pržulj N, Corneil DG, Jurisica I. Modeling interactome: scale-free or geometric? *Bioinformatics*. 2004; 20(18):3508–3515. <https://doi.org/10.1093/bioinformatics/bth436> PMID: 15284103
9. Pržulj N. Biological network comparison using graphlet degree distribution. *Bioinformatics*. 2007; 23(2): e177–e183. <https://doi.org/10.1093/bioinformatics/btl301> PMID: 17237089
10. Yaveroğlu ÖN, Malod-Dognin N, Davis D, Levnajic Z, Janjic V, Karapandza R, et al. Revealing the hidden language of complex networks. *Scientific reports*. 2014; 4:4547. <https://doi.org/10.1038/srep04547> PMID: 24686408
11. Kuchaiev O, Milenković T, Memišević V, Hayes W, Pržulj N. Topological network alignment uncovers biological function and phylogeny. *Journal of The Royal Society Interface*. 2010; 7(50):1341–1354. <https://doi.org/10.1098/rsif.2010.0063>
12. Malod-Dognin N, Pržulj N. L-GRAAL: Lagrangian Graphlet-based Network Aligner. *Bioinformatics*. 2015; <https://doi.org/10.1093/bioinformatics/btv130>
13. Saraph V, Milenković T. MAGNA: maximizing accuracy in global network alignment. *Bioinformatics*. 2014; 30(20):2931–2940. <https://doi.org/10.1093/bioinformatics/btu409> PMID: 25015987
14. Mamano N, Hayes W. SANA: Simulated Annealing far outperforms many other search algorithms for biological network alignment. *Bioinformatics*. 2017; 0(0):8.
15. Hočevar T, Demšar J. A combinatorial approach to graphlet counting. *Bioinformatics*. 2014; 30(4): 559–565. <https://doi.org/10.1093/bioinformatics/btt717> PMID: 24336411
16. Melckenbeeck I, Audenaert P, Michoel T, Colle D, Pickavet M. An Algorithm to Automatically Generate the Combinatorial Orbit Counting Equations. *PLoS ONE*. 2016; 11(1). <https://doi.org/10.1371/journal.pone.0147078> PMID: 26797021
17. Chatr-aryamontri A, Breitkreutz BJ, Heinicke S, Boucher L, Winter A, Stark C, et al. The BioGRID interaction database: 2013 update. *Nucleic Acids Research*. 2013; 41(D1):D816–D823. <https://doi.org/10.1093/nar/gks1158> PMID: 23203989
18. Pillich RT, Chen J, Rynkov V, Welker D, Pratt D. NDEx: A Community Resource for Sharing and Publishing of Biological Networks. *Protein Bioinformatics: From Protein Modifications and Networks to Proteomics*. 2017; p. 271–301. [https://doi.org/10.1007/978-1-4939-6783-4\\_13](https://doi.org/10.1007/978-1-4939-6783-4_13)

19. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos JS, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics*. 2009; 10:421. <https://doi.org/10.1186/1471-2105-10-421> PMID: 20003500
20. Rahman M, Bhuiyan MA, Al Hasan M. Graft: An efficient graphlet counting method for large graph analysis. *IEEE Transactions on Knowledge and Data Engineering*. 2014; 26(10):2466–2478. <https://doi.org/10.1109/TKDE.2013.2297929>
21. Pržulj N, Corneil DG, Jurisica I. Efficient estimation of graphlet frequency distributions in protein—protein interaction networks. *Bioinformatics*. 2006; 22(8):974–980. <https://doi.org/10.1093/bioinformatics/btl030> PMID: 16452112
22. Sloane N. Online Encyclopedia of Integer Sequences (OEIS);. Available from: <http://oeis.org/A000088>.
23. Mckay BD. Nauty; 2010. Available from: <http://users.cecs.anu.edu.au/~bdm/nauty>.
24. Codenotti P, Katebi H, Sakallah KA, Markov IL. Conflict Analysis and Branching Heuristics in the Search for Graph Automorphisms. In: *Tools with Artificial Intelligence (ICTAI)*. IEEE; 2013.
25. Gross JL. Graph Theory—Lecture 2: Structure and Representation—Part A; Available from: <http://www.cs.columbia.edu/~cs4203/files/GT-Lec2.pdf>.
26. Automorphism of a group;. Available from: [https://groupprops.subwiki.org/wiki/Automorphism\\_of\\_a\\_group](https://groupprops.subwiki.org/wiki/Automorphism_of_a_group).
27. Sloane N. Online Encyclopedia of Integer Sequences (OEIS);. Available from: <http://oeis.org/A000666>.