

# Triadic analysis of affiliation networks

JASON CORY BRUNSON

Center for Quantitative Medicine, UConn Health, Farmington, CT 06030

(*e-mail*: brunson@uchc.edu)

---

## Abstract

Triadic closure has been conceptualized and measured in a variety of ways, most famously the clustering coefficient. Existing extensions to affiliation networks, however, are sensitive to repeat group attendance, which manifests in bipartite models as biclique proliferation. Whereas this sensitivity does not reflect common interpretations of triadic closure in social networks, this paper proposes a measure of triadic closure in affiliation networks designed to control for it. To avoid arbitrariness, the paper introduces a triadic framework for affiliation networks, within which a range of measures can be defined; it then presents a set of basic axioms that suffice to narrow this range to the one measure. An instrumental assessment compares the proposed and two existing measures for reliability, validity, redundancy, and practicality. All three measures then take part in an investigation of three empirical social networks, which illustrates their differences.

---

## 1 Introduction

Triadic analysis, which emphasizes the interactions within subsets of three nodes, has long been central to network science. Meanwhile, affiliation (or co-occurrence) data have long been a major source of empirical networks. Most triadic analyses of affiliation networks either collapse their higher-order structure or focus on relations among triples of nodes, often of mixed type. This paper, building upon some recent contributions, focuses instead on triples of actors, together with the non-actor structure that establishes relations among them.

### 1.1 Background

**Precursors** Previous triadic approaches in the social networks literature provide examples of hypothesis formulation, measure design, and sociological interpretation that inspired the present analysis. One thread begins with a series of studies designed to test socio-structural predictions of cognitive balance theory (Davis, 1967). These predictions apply at the level of triads, but could be analyzed statistically by aggregating over an entire graph. For example, the transitive property, under which the directed relations  $p \rightarrow q \rightarrow r$  imply the relation  $p \rightarrow r$ , describes social graphs with a specific hierarchical structure (Holland & Leinhardt, 1971). While this structure would be hard to measure directly, the transitivity ratio (the global proportion of instances of  $p \rightarrow q \rightarrow r$  in which  $p \rightarrow r$ ) provides a simple measure of how closely a graph respects this property (Harary & Kommel, 1979).

A separate thread concerns the “small world” property, a high concentration of ties within communities yet counterintuitively low distances between actors in different communities, observed in empirical social networks (de Sola Pool & Kochen, 1978). The “strong triadic closure” (STC) hypothesis proposed to reconcile these properties by ascribing a cohesion role to strong ties within communities and a bridging role to weak ties between them (Granovetter, 1973). STC distinguishes two levels of tie (strong and weak) and posits that strong ties lead to more closures. A reorientation from triads to ego networks led to the “structural holes” framework, in which an actor with many weak ties, hence a more disconnected neighborhood, has increased potential as a broker. The local measure of constraint was introduced to quantify how these neighborhood connections limit brokerage potential (Burt, 1992). A later, independent study introduced the similar but simpler clustering coefficient to quantify “cliquishness” across a family of small world models (Watts & Strogatz, 1998).

**Conventions** The present study concerns social networks, but the concepts generalize to any affiliation network (AN) setting. Most terminology and notation is taken from standard references (Bondy & Murty, 1976; Wasserman & Faust, 1994). Additional concepts will be defined as needed.

A graph  $G = (V, E)$  consists of a finite set  $V$  of nodes and a set  $E \subseteq V \times V$  of edges  $e = (v, w)$ . Edges will be symmetric and will not include duplicates or loops. A graph is bipartite if its nodes can be partitioned into subsets  $V_1$  and  $V_2$  in such a way that  $E \subseteq V_1 \times V_2$ . The degree of a node  $v$  is the number of edges containing  $v$ . A subgraph of  $G$  is a graph  $G' = (V', E')$  satisfying  $V' \subseteq V$  and  $E' \subseteq E$ , and a subset  $W \subseteq V$  of nodes induces the subgraph  $(W, E \cap (W \times W))$ .

Traditional social networks consist of actors having (here, symmetric) relations among them, and are modeled as graphs with actors represented by nodes and relations by edges. Three actors, together with the relations among them, form a triad. The triads of a traditional network  $G$  take four types  $i = 0, 1, 2, 3$ , according to the number of relations among their actors; the tallies  $s_i = s_i(G)$  of each type constitute the *triad census*  $(s_0, s_1, s_2, s_3)$ . The (*classical*) *clustering coefficient*, often described as the proportion of connected triples that are closed (Newman, 2003), is then the ratio  $C(G) = 3 \times s_3 / (s_2 + 3 \times s_3)$ .

Relations among the actors of an AN are established through common attendance at events; each event is attended by some subset of actors. ANs are modeled as bipartite graphs,  $V_1$  consisting of the actors and  $V_2$  the events. Though actors are only tied to events, in both settings the neighbors of an actor  $v$  shall be the actors related to  $v$ . If actors who coattended events are assigned edges, then they (without the events) form a traditional social network called the projection.

**Organization** Sec. 1.2 proposes the new clustering coefficient. The main body of the paper is split between theoretical (Sec. 2) and empirical (Sec. 3) assessments of this statistic, and begin with their own organizational outlines. In short, Sec. 2 explores triadic analysis in the abstract, while Sec. 3 performs triadic analyses on empirical data. Sec. 4 summarizes the paper, its limitations, and future directions. All analyses are performed, and images produced, using the open-source statistical programming language R, with the `igraph` and `ggplot2` packages in particular (R Development Core Team, 2008; Csardi & Nepusz,

2006; Wickham, 2009). Full code is available at <https://github.com/corybrunson/triadic>.

## 1.2 The exclusive clustering coefficient

**Motivation** “Triadic closure” (TC) refers to the tendency for the relations  $(p, q)$  and  $(q, r)$  to entail the relation  $(p, r)$ . This entailment need not be causal, nor even chronological, but interpersonal interpretations of TC posit that the common neighbor  $q$  facilitates, or even initiates, the connection between  $p$  and  $r$ . Such interpretations, however, are at odds with common measures of TC, especially in the AN setting.

The clustering coefficient, for example, is often evaluated on projections; this shall be the meaning of the shorthand  $C(G)$  when  $G$  is an AN.<sup>1</sup> A conspicuous feature of these projections is the proliferation of clique graphs  $K_n$ , which consist of  $n$  nodes and all  $\binom{n}{2}$  possible edges between them.  $n$  actors in  $G$  who attend any single event produce a copy of  $K_n$  in the projection, which contains  $\binom{n}{3}$  3-edge triads. These can dramatically increase  $C(G)$ , so that its values are often largely determined by event size (Newman, 2001; Glänzel & Schubert, 2004). High event attendance, however, does not guarantee TC: Individuals in distinct, pre-existing social groups at a common event may interact primarily with others in their own groups, and forge few if any inter-group relations.

Attempts to account for this inflation of  $C$  have taken both “conversion” (at the projection level) and “direct” (at the AN level) approaches. Conversion approaches have, for example, standardized the value of  $C$  by its values at a suitable null model (Uzzi & Spiro, 2005) and applied clustering coefficients designed for weighted networks to weighted projections (Saramäki *et al.*, 2007). These methods help discriminate levels of TC among ANs, but at some cost to interpretability.

Two recent direct approaches define new clustering coefficients in terms of AN structure among triples of actors (Opsahl, 2013; Liebig & Rao, 2014). The *Opsahl clustering coefficient*  $C^*$ , for example, restricts the notion of “connected triples” (of actors) to those who are pairwise connected through distinct events. It can be defined as the proportion of 4-paths that are closed: The graph  $P_d$  consisting of distinct nodes  $v_0, v_1, \dots, v_d$  and edges  $(v_i, v_{i+1})$  is called the  $d$ -path; if, instead,  $v_0 = v_d$ , the result is the  $d$ -cycle  $C_d$ .<sup>2</sup> (Both have  $d$  edges; see Fig. 1c,d.) For a 4-path in  $G$  to be “closed” means for it to be contained in a 6-cycle.<sup>3</sup> In an empirical test,  $C^*$  took much smaller values than  $C$ , and the two statistics diverged most on the network with the greatest mean event size (Opsahl, 2013).

However, these measures may still be at odds with the popular interpretation of TC: The same pre-existing groups that attend one event are likely to attend others, though this no more entails TC than attendance at the first. Such repeat group attendance manifests in bipartite AN models as the proliferation of biclique graphs  $K_{n,m}$ , which consist of  $n$  actors who each attend each of  $m$  events (hence  $n \times m$  edges). Indeed, bicliques and similar motifs have been observed in empirical ANs at frequencies greater than expected by chance

<sup>1</sup>  $C$  evaluates to zero on any bipartite graph.

<sup>2</sup> The  $d$ -paths involved in this calculation must begin and end at actor nodes.

<sup>3</sup> Several other studies have proposed bipartite clustering coefficients that concern triples of nodes but not of actors, and are not considered here (Opsahl, 2013).

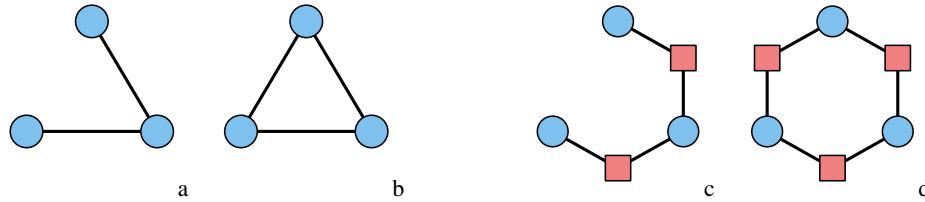


Fig. 1: Traditional and affiliation network conceptions of triadic closure: In traditional networks, the 2-edge triad (a) is “open” while the 3-edge triad (b) is “closed”. The clustering coefficient  $C$  is defined either as a ratio of the numbers of triads of these types (see Sec. 1.1) or as the proportion of subgraphs of the form (a) that are contained in a subgraph of the form (b). One extension of this idea to ANs uses the 4-path (c) and the 6-cycle (d) in place of these triads. The Opsahl clustering coefficient  $C^*$  is defined as the proportion of subgraphs (c) that are contained in a subgraph (d). (Circular nodes denote actors; square nodes denote events.)

(Borgatti & Everett, 1997; Carrino, 2006). Just as  $C$  is sensitive to cliques,  $C^*$  is sensitive to bicliques: If  $m \geq 3$ , then  $K_{3,m}$  contains  $6m(m-1)$  4-paths, each of which is closed; the effect grows geometrically with  $n$ .<sup>4</sup> Thus empirical values of  $C^*$  may be dominated by patterns of repeat group attendance. The need for a measure of TC in ANs that also controls for this artifact motivated the present study.

**Proposal** The proposed graph statistic follows  $C^*$  in restricting to pairwise connectivity through separate events within a triad. It also addresses two concerns raised by  $C^*$ : First, the 4-paths and 6-cycles in its calculation contain no intermediate edges—each event is attended by only two of the three actors.<sup>5</sup> This eliminates the direct influence of repeat group attendance.

Another concern is how the population of actors (or of triads) should be weighted in the calculation.  $C$  weights all actors equally, in that any ordered triple of actors can have at most one 2-path through them in the projection. In contrast, many 4-paths may exist through a single ordered triple in an AN, due to a multiplicity of shared events, so that more prolific actors will tend to have more influence on the value of  $C^*$ . Because the present study takes an actor-centric approach, the proposed statistic is designed to weight actors equally.

The statistic is denoted  $C^\circ$ . It asks, *provided  $p$  and  $q$  attend some event without  $r$ , and  $q$  and  $r$  attend some event without  $p$ , what is the probability that  $p$  and  $r$  attend some event without  $q$ ?* Since  $C^\circ$  measures TC only through pairwise-exclusive events, it shall be called the *exclusive clustering coefficient*.

#### Example 1.1

Fig. 2 depicts four ANs that project to the same the “kite” graph. AN (a) exhibits TC in the sense of interest, while (b) exhibits TC of the kind  $C^*$  was designed to ignore.  $C^*$  evaluates to  $\frac{3}{5}$  at (a) and to 0 at (b), and  $C^\circ$  agrees on both.  $C^*$  evaluates to  $\frac{5}{8}$  at (c), and to  $\frac{3}{4}$  at (d),

<sup>4</sup> Though note that  $K_{3,2}$  contains  $6 \times 2 = 12$  open (and no closed) 4-paths.

<sup>5</sup> This choice, and some alternatives, have received their own treatment (Liebig & Rao, 2014).

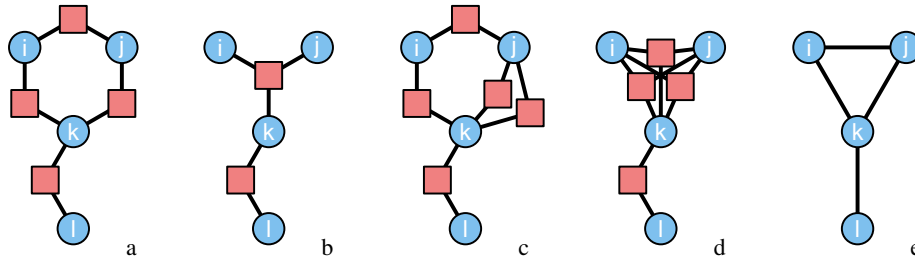


Fig. 2: Four ANs (a–d) having the same projection (e).

due to additional copies of  $P_4$  and  $C_6$ . For instance, six copies of  $P_4$  in (d) proceed from  $i$  through  $j$  to  $k$ , and each is closed. In contrast,  $C^\circ$  takes the familiar values  $\frac{3}{5}$  at (c) and 0 at (d), since it is calculated on the same numbers of distinct 4-paths and 6-cycles.

By eliminating sources of 3-edge triads other than the popular meaning of TC,  $C^\circ$  may help to infer dynamic information from static data. The popular meaning is dynamic: Actors who are not neighbors, but who have neighbors in common at one time, become neighbors at a later time (Easley & Kleinberg, 2010; Martin *et al.*, 2013). If a traditional network  $G$  has edges labeled by instants in time, such that an edge labeled  $t$  is said to exist after  $t$  but not before, define the *dynamic triadic closure*  $D(G)$  to be, among those triads at which there is at some time an open 2-path, the proportion at which there is at a later time a 3-cycle. If  $G$  is an AN with events labeled by time, then  $D$  shall be calculated on its projection, where each edge is labeled by the earliest event that projects to it.

In the traditional setting, if a network has time-labeled edges, no two of which are simultaneous, then  $D = s_3 / (s_2 + s_3) = C / (3 - 2 \times C)$ . In the AN setting, pairwise-exclusive events are essential to  $D$ , since an open 2-path in the projection must correspond to a triad with only pairwise-exclusive events. While the two calculations are in general unequal even when no two events are simultaneous,  $C^\circ$  could provide a useful estimate of  $D$ .

## 2 Theoretical analyses

This section formalizes the exclusive clustering coefficient and evaluates its theoretical merits. Sec. 2.1 develops a formal notion of “triad” for ANs. On this foundation, Sec. 2.2 unifies  $C$ ,  $C^*$ , and  $C^\circ$  into a generic clustering coefficient. This definition specializes to impracticably many statistics, which Sec. 2.3 whittles down by appeal to several properties suited to present purposes. The technical details of this process are relegated to Sec. 2.4.

### 2.1 Triads

**Scheduled subgraphs** A triad-centric approach to ANs requires an object of study. What, then, is a “triad”? Since the triads of a traditional network are the subgraphs *induced* by three actors, a suitable analog of induced subgraphs for ANs would suffice. This paper proposes to include those events that establish relations among a set of actors:

*Definition 2.1*

Given an AN  $G$  and a subset  $W$  of actors of  $G$ , the subgraph of  $G$  *scheduled* by  $W$  is the subgraph induced by the actors  $W$  together with all events attended by at least two actors in  $W$ . *Scheduled* graph maps are defined analogously to induced graph maps, and the *triads* of  $G$  are the subgraphs scheduled by three actors.

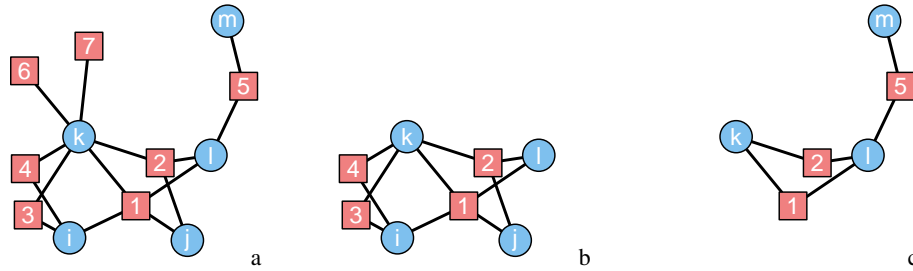


Fig. 3: An AN (a) of actors  $i, j, k, l, m$  and its scheduled subgraphs (b) at  $\{i, j, k, l\}$  and (c) at  $\{k, l, m\}$ .

### Example 2.2

Fig. 3 depicts an AN of five actors and two of its scheduled subgraphs. Note in particular that the scheduled subgraph on the entire set of actors (not shown) does not include events 6 and 7, since they play no role in establishing relations among the actors.<sup>6</sup>

**Triad censuses** The classification of AN triads is straightforward but not trivial. While traditional triads fall into four isomorphism classes (see Sec. 2.2), AN triads, in theory, occupy arbitrarily many, due to the unlimited number of events two or three actors might attend. Consider an arbitrary triad with actors  $p, q, r$ . Take  $w_{pq}$  to be the number of events attended by  $p$  and  $q$ , similarly define  $w_{qr}$  and  $w_{pr}$ , and take  $w_{pqr}$  to be the number of events attended by all three. (Note that  $w_{pq}$  does not depend on  $r$ , etc.) Up to isomorphism, it may be assumed that  $w_{pq} \geq w_{qr} \geq w_{pr}$  (otherwise relabel the actors). Necessarily,  $w_{pqr} \leq w_{pq}$ . Let  $\mu = (\mu_1, \mu_2, \mu_3) = (w_{pq} - w_{pqr}, w_{qr} - w_{pqr}, w_{pr} - w_{pqr})$  count the “exclusive” events between each pair of actors, and let  $w = w_{pqr}$  count the “inclusive” events attended by all three. The pair  $(\mu, w)$  determines the isomorphism class of the triad.<sup>7</sup> Since  $\mu_1 \geq \mu_2 \geq \mu_3$ ,  $\mu$  is an integer partition of three parts; write  $\mu \in \text{Par}_3$ . Where  $\mathbb{Z}_{\geq 0}$  is the set of nonnegative integers and  $\mathcal{T}$  is the set of triad isomorphism classes, this gives a bijective correspondence

$$\mathcal{T} \leftrightarrow \text{Par}_3 \times \mathbb{Z}_{\geq 0}.$$

Write  $\text{Tr}_{\mu w}$  for the triad described above, and  $s_{\mu w} = s_{\mu w}(G)$  for the number of triads of  $G$  isomorphic to  $\text{Tr}_{\mu w}$ . The (full) triad census of  $G$  is then the array  $(s_{\mu w})_{\mu, w}$ . The partitions  $\text{Par}_3$  can be totally ordered, and thereby the census arranged in a matrix, whose

<sup>6</sup> This shows that an AN need not be the scheduled subgraph of its actors, contrary to the analogous property of induced subgraphs. Their projections, however, are the same (up to edge weights).

<sup>7</sup> While this scheme is more intuitive, a more storage-friendly enumeration of the triad classes is given by the quadruple  $(w_{pqr} - w_{pq}, w_{pq} - w_{qr}, w_{qr} - w_{pr}, w_{pr}) \in (\mathbb{Z}_{\geq 0})^4$ .

size depends on the network.<sup>8</sup> Necessarily,  $\sum_{\mu,w} s_{\mu w} = \binom{|V_1|}{3}$ . The triads scheduled from  $i, j, k$  in Fig. 2 (a–d), for example, are  $\text{Tr}_{(1,1,1),0}$ ,  $\text{Tr}_{(0,0,0),1}$ ,  $\text{Tr}_{(2,1,1),0}$ , and  $\text{Tr}_{(0,0,0),3}$ .

This scheme explodes as networks grow dense. The following alternative scheme is instead bounded, but nonetheless captures useful affiliation structure: The events of a triad fall into four structural equivalence classes, according to which actors attend them. Instead of binning triads by *how many* events they have in each class, bin them by whether they contain *some* event in each class. If  $\text{Tr}_{\mu w}$  has, for example, any inclusive event (i.e., if  $w > 0$ ), then  $\text{Tr}_{\mu w}$  shares a bin with  $\text{Tr}_{\mu,1}$ ; otherwise it is  $\text{Tr}_{\mu,0}$ . Each bin then contains exactly one representative  $\text{Tr}_{\mu w}$  with  $\mu_1, \mu_2, \mu_3, w \in \{0, 1\}$ , and this bin is determined by the two numbers  $x = \mu_1 + \mu_2 + \mu_3 \in \{0, 1, 2, 3\}$  and  $y = w \in \{0, 1\}$ . The *structural triad census* consists of the eight tallies  $t_{xy}$  of triads in each bin. Though containing only twice as many bins as the simple census, the structural census contains useful additional information (see Thm. 2.7 and Sec. 3.1).

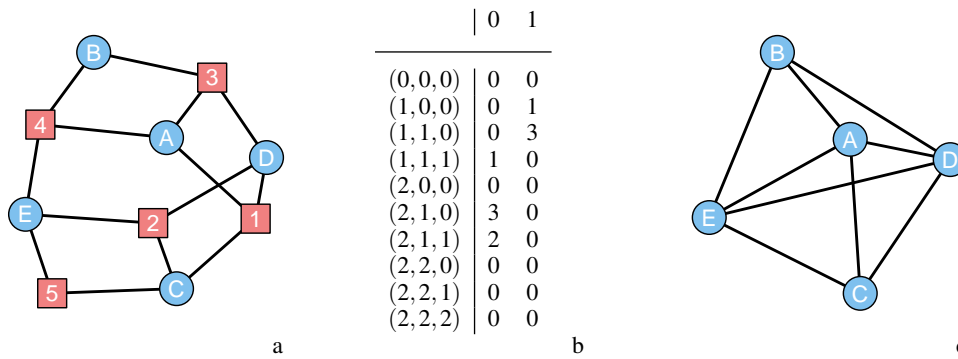


Fig. 4: The network DG2 (a), its full triad census (b), and its projection (c). The column in (b) indicates the number of inclusive events; the row indicates the distribution of exclusive events across pairs of actors. For example, the triad at (A, B, C) is tallied in column 0, row (2, 1, 0) (see Ex. 2.3).

*Example 2.3*

The network DG2, depicted in Fig. 4 with its full census and its projection, is taken from a famous study of the American racial caste system (Davis *et al.*, 1941). As an example of a social unit, the study presented attendance data for five acquainted women (“Miss A” through “Miss E”) and five social activities (bridge, dinner, movies, dance, and visiting), forming an AN. The projection contains three 2-edge and seven 3-edge triads, so the simple census is (0, 0, 3, 7). (Therefore, incidentally,  $C(\text{DG2}) = \frac{3 \times 7}{3 + 3 \times 7} = \frac{7}{8}$ .)

These tallies obscure higher-order structure: The seven fully-connected triads fall into four classes. One might be called “symmetric” and “exclusive”: Mss. B, D, and E attended

<sup>8</sup> Where  $n = \max(\mu)$ , there are bijections  $\sigma : \text{Par}_3^{(n)} \rightarrow \left\{ \binom{n+3}{3} \right\}$ , from the partitions in  $\text{Par}_3$  having parts  $\leq n$  to the subsets of  $\{1, \dots, n+3\}$  of size 3, and  $\rho : \left\{ \binom{n+3}{3} \right\} \rightarrow \{1, \dots, \binom{n+3}{3}\}$ , which indexes these subsets; the composition  $\rho \circ \sigma : \text{Par}_3^{(n)} \rightarrow \{1, \dots, \binom{n+3}{3}\}$  indexes the partitions.  $\sigma$  is a classical bijection (Stanley, 2002);  $\rho$  is the *revolving door ordering* (Kreher & Stinson, 1999).

no events together, but each pair were present at one, so that their triad is (isomorphic to)  $\text{Tr}_{(1,1,1),0}$ . Two triads were exclusive but not symmetric: Mss. C and E attended two events without Ms. A, though Ms. A attended different, separate events with Mss. C and E; they thus form a triad  $\text{Tr}_{(2,1,1),0}$ , as do Mss. A, D, and E. The remaining four were “inclusive”, in that all three women attended some event together (specifically, the four activities of attendance 3). In each case, at least one pair of women attended another event together, forming the triads  $\text{Tr}_{(1,0,0),1}$  and  $\text{Tr}_{(1,1,0),1}$ . The women constituting each of the three 2-edge triads joined in no single activity together, instead forming three copies of  $\text{Tr}_{(2,1,0),0}$ . For example, Mss. A and B attended two events together, Mss. A and C one, and Mss. B and C none. (As an exercise, the reader might recover the structural census from Fig. 4a,b.)

## 2.2 Category framework

**Graph maps** A generic clustering coefficient will be defined in terms of graph maps. For present purposes, a graph map  $\phi : G \rightarrow H$  (“ $\phi$  from  $G$  to  $H$ ”) shall assign each node  $v$  of  $G$  to a node  $\phi(v)$  in  $H$  (the image of  $v$  under  $\phi$ ) in such a way that every edge  $(v, w)$  in  $G$  is preserved, i.e.  $(\phi(v), \phi(w))$  is an edge in  $H$ . One example is the inclusion of a subgraph  $G \subseteq H$ . A graph map  $\phi : G \rightarrow H$  is called induced if the image  $\phi(G) \subseteq H$  is an induced subgraph. The images  $\phi(v)$  in  $H$  and the preserved edges among them form the image of  $G$  in  $H$ . Two graph maps  $\phi : G \rightarrow H$  and  $\psi : H \rightarrow K$  yield the composition  $\psi \circ \phi : G \rightarrow K$  defined by  $\psi \circ \phi(v) = \psi(\phi(v))$ . Such a graph map  $\psi \circ \phi : G \rightarrow K$  is said to factor through  $H$ ; for example, any map  $\phi : G \rightarrow H$  factors through its image  $\phi(G) \subseteq H$ .

A graph map  $\phi : G \rightarrow H$  is injective if it sends no two nodes in  $G$  to the same node in  $H$ , and surjective if every node in  $H$  is the image of some node in  $G$  (its pre-image); by a “copy” of  $G$  in  $H$ , or a path or cycle “in  $G$ ”, shall be meant the image of an injective map. (By convention, paths and cycles in an AN arise from maps that send  $v_0$  to an actor.) Thus a 4-path  $\phi : P_4 \rightarrow G$  is closed if it factors through  $C_6$ .

An injective, surjective map is bijective, and a bijective map  $\phi : G \rightarrow H$  is an isomorphism if it is induced—that is, if it preserves absences of edges ( $(\phi(v), \phi(w)) \notin G$  whenever  $(v, w) \notin H$ ). The isomorphisms establish an equivalence relation on graphs; two graphs related by an isomorphism are said to be isomorphic, and to lie in the same isomorphism class. Two nodes  $v, w \in G$  are structurally equivalent if there is an isomorphism  $G \rightarrow G$  that sends every node to itself except exchanges  $v$  and  $w$ ; this establishes an equivalence relation on the nodes of  $G$ .

**Categories** The framework of category theory, though not necessary, absorbs some useful and unobjectionable yet messy assumptions into the notation, provides a catalogue of natural examples, and avoids unnecessary constraints on the range of possibilities.<sup>9</sup>

<sup>9</sup> While there are infinitely many AN triads, their combinatorial complexity is limited (see Sec. 2.1). It would be short work to classify a useful collection (19, by the author’s count) of clustering coefficients, in the sense of Def. 2.4 and including  $C$ ,  $C^*$ , and  $C^\circ$ , by which structural equivalence classes of events the events of  $W$  and  $X$  may be mapped to, and which of these should then be considered congruent. This scheme, however, would omit more ad hoc clustering coefficients, for instance one that requires the events  $v_1, v_3$  of  $W$  to be mapped to exclusive events but places no



A category  $\mathcal{C}$  consists of a set of objects; for each pair of objects  $A, B$ , a set  $\text{Hom}_{\mathcal{C}}(A, B)$  of morphisms from  $A$  to  $B$ ; and, for each pair of morphisms  $f \in \text{Hom}_{\mathcal{C}}(B, C)$  and  $g \in \text{Hom}_{\mathcal{C}}(A, B)$ , the composition  $f \circ g \in \text{Hom}_{\mathcal{C}}(A, C)$ ; all subject to the following conditions (Mitchell, 1965):

- i. (Identity) For every  $A \in \mathcal{C}$  there exists  $\text{id}_A \in \text{Hom}_{\mathcal{C}}(A, A)$  satisfying  $f \circ \text{id}_A = f$  and  $\text{id}_A \circ g = g$  for any  $f \in \text{Hom}_{\mathcal{C}}(A, B)$  or  $g \in \text{Hom}_{\mathcal{C}}(C, A)$ .
- ii. (Associativity) For any triple of morphisms  $f \in \text{Hom}_{\mathcal{C}}(C, D)$ ,  $g \in \text{Hom}_{\mathcal{C}}(B, C)$ , and  $h \in \text{Hom}_{\mathcal{C}}(A, B)$ ,  $f \circ (g \circ h) = (f \circ g) \circ h$ .

A subcategory  $\mathcal{C}' \subseteq \mathcal{C}$  consists of the same objects as  $\mathcal{C}$  and subsets  $\text{Hom}_{\mathcal{C}'}(A, B) \subseteq \text{Hom}_{\mathcal{C}}(A, B)$  that also form a category. A congruence relation  $\sim$  on  $\mathcal{C}$  consists of equivalence relations  $\sim_{A, B}$  on each  $\text{Hom}_{\mathcal{C}}(A, B)$  that are compatible with the composition of morphisms, so that the quotient category  $\mathcal{C}/\sim$  is determined by the objects of  $\mathcal{C}$  and the equivalence classes of morphisms of  $\mathcal{C}$  under  $\sim$ .

Henceforth, view  $\mathcal{T}$  as the category of AN triads, with morphisms the graph maps  $\phi : H \rightarrow K$  that assign the actors of  $H$  to distinct actors of  $K$  (and therefore send events only to events), and with composition given by  $(f \circ g)(v) = f(g(v))$ .  $\mathcal{T}$  can be viewed as a subcategory of the category of graphs (Hell, 1979) (with many objects omitted). Write  $\text{Hom}_{\mathcal{T}}^K(G, H)$  for the set of morphisms from  $G$  to  $H$  that factor through  $K$ . If  $G$  is any AN, write  $\text{Hom}_{\mathcal{T}}(H, G)$  (an abuse of notation) for the set of all morphisms from  $H$  to any triad of  $G$ .

**Clustering coefficients** All three clustering coefficients described in Sec. 1.2 are expressible in category-theoretic terms. Let  $\approx$  denote the congruence relation on  $\mathcal{T}$  given by taking any two maps that agree on actors to be congruent. For example, there is only one graph map from  $P_4$  to the kite graph (a) in Fig. 2 that sends  $v_0, v_2, v_4$  to  $i, j, k$  (respectively), and likewise only one such map to (b). However, there are several such maps to (c), which are all congruent in  $\mathcal{T}/\approx$ . Thus  $\approx$  is a “strong” relation in that it relates very many morphisms. It turns out that, for an affiliation network  $G$ ,

$$C(G) = \frac{|\text{Hom}_{\mathcal{T}/\approx}^{C_6}(P_4, G)|}{|\text{Hom}_{\mathcal{T}/\approx}(P_4, G)|} = \frac{|\text{Hom}_{\mathcal{T}/\approx}(C_6, G)|}{|\text{Hom}_{\mathcal{T}/\approx}(P_4, G)|}. \quad (1)$$

The Opsahl clustering coefficient restricts the morphisms in Eq. 1 to injective graph maps. It is straightforward to check that these form a subcategory  $\widetilde{\mathcal{T}} \subset \mathcal{T}$ . No congruence relation was imposed; for consistency of notation, write  $\widetilde{\mathcal{T}}/\approx$  for  $\widetilde{\mathcal{T}}$ , where  $\approx$  denotes equality of graph maps (the weakest possible relation).  $C^*$  is then realized as

$$C^*(G) = \frac{|\text{Hom}_{\widetilde{\mathcal{T}}/\approx}^{C_6}(P_4, G)|}{|\text{Hom}_{\widetilde{\mathcal{T}}/\approx}(P_4, G)|}, \quad (2)$$

analogously to the first formulation in Eq. 1. The present proposal further restricts the morphisms in Eq. 2 to induced injections. These turn out to form their own subcategory

such constraint on  $v_5$  in  $X$ . Such a statistic would violate Axiom 1, but may be very useful in certain settings (compare to the discussion of STC in Sec. 3.2).

$\overline{\mathcal{T}} \subset \widetilde{\mathcal{T}}$ . Additionally, the graph maps that agree on actors *and that send events to structurally equivalent images* constitute a congruence relation  $\simeq$  on  $\overline{\mathcal{T}}$  (or  $\mathcal{T}$ ), which is weaker than  $\approx$  but stronger than  $=$ . The statistic  $C^\circ$  is then realized as

$$C^\circ(G) = \frac{|\text{Hom}_{\overline{\mathcal{T}}/\simeq}^{C_6}(P_4, G)|}{|\text{Hom}_{\overline{\mathcal{T}}/\simeq}(P_4, G)|} = \frac{|\text{Hom}_{\overline{\mathcal{T}}/\simeq}(C_6, G)|}{|\text{Hom}_{\overline{\mathcal{T}}/\simeq}(P_4, G)|}. \quad (3)$$

### 2.3 Axiomatic approach

**General formulation** What is a “clustering coefficient”, especially in the AN setting? Sec. 2.2 formulated three variations on the idea, and this section presents a single unifying definition.

The statistics  $C$  and  $C^*$  differ in three respects: the choice between the formulations in Eq. 1 (which sometimes agree), the subcategory of graph maps from which the morphisms in Eq. 1 are drawn, and the congruence relation imposed on them. Whereas  $P_4$  (isomorphic to  $\text{Tr}_{(1,1,0),0}$ ) and  $C_6$  (isomorphic to  $\text{Tr}_{(1,1,1),0}$ ) are now recognizable as two among an infinite collection of triads (see Fig. 5), a fourth choice presents: What makes a triple of actors “open” or “closed”? Another direct approach (Liebig & Rao, 2014) considered three alternatives to  $C_6$ :  $\text{Tr}_{(1,1,0),1}$ ,  $\text{Tr}_{(1,0,0),2}$ , and  $\text{Tr}_{(0,0,0),3}$ . (These are the four AN triads whose duals are also triads, and in fact are self-dual (Breiger, 1974).) Alternatives to  $P_4$ , sometimes taken in pairs, were obtained by removing a single event from these. The four choices thus outlined are incorporated into the following general definition:

#### Definition 2.4

Pick canonical triads  $X \in \mathcal{T}$  and  $W \subset X$ , a canonical subgraph relation  $\iota : W \rightarrow X$  (there may be many), a subcategory  $\mathcal{C} \subseteq \mathcal{T}$ , and a congruence relation  $\sim$  on  $\mathcal{C}$ . A (*global*) *clustering coefficient* of  $G$  shall be a statistic of either form

$$\widehat{C}(G) = \frac{|\text{Hom}_{\mathcal{C}/\sim}^X(W, G)|}{|\text{Hom}_{\mathcal{C}/\sim}(W, G)|} \quad (\text{“rate of wedge closure”}) \quad \text{or} \quad (4)$$

$$\widehat{C}(G) = \frac{|\text{Hom}_{\mathcal{C}/\sim}(X, G)|}{|\text{Hom}_{\mathcal{C}/\sim}(W, G)|} \quad (\text{“alcove-to-wedge ratio”}), \quad (5)$$

where morphisms factor through  $X$  only via  $\iota$ . Call the morphisms  $\text{Hom}_{\mathcal{C}/\sim}(W, G)$  the *wedges* of  $G$ —*closed* if they factor through  $X$ , *open* if not—and  $\text{Hom}_{\mathcal{C}/\sim}(X, G)$  the *alcoves* of  $G$ .

Further designate a *center* actor  $v_c \in \{p, q, r\}$  in (each)  $W$ . Given an actor  $j \in G$ , obtain the (*local*) *clustering coefficient*  $\widehat{C}(j)$  of  $j$  by requiring of the morphisms in Eq. 4 or 5 that  $\phi(v_c) = j$  and  $\psi(\iota(v_c)) = j$ —that is, that wedges and alcoves are *centered* at  $j$ . The *wedge-dependent clustering coefficient*  $\widehat{C}_\ell$  of an affiliation network  $G$  shall be the mean value of  $\widehat{C}(j)$  across the actors  $j$  at which exactly  $\ell$  wedges are centered.

By letting  $X$  range over the four self-dual triads;  $\mathcal{C}$  over  $\mathcal{T} \supseteq \widetilde{\mathcal{T}} \supseteq \overline{\mathcal{T}}$ ;  $\sim$  over  $=$ ,  $\simeq$ , and  $\approx$ ; and adopting either Eq. 4 or 5, Def. 2.4 specializes to  $4 \times 3 \times 3 \times 2 = 72$

fairly straightforward statistics, which include  $C$ ,  $C^*$ , and  $C^\circ$ .<sup>10</sup> For present purposes, the best choice of  $X$  is clearly  $\text{Tr}_{(1,1,1),0}$ , leaving  $W = \text{Tr}_{(1,1,0),0}$ . These choices are assumed henceforth. (Note, however, that Thm. 2.6 does not require this assumption.)

Table 1: Three measures of global and local triadic closure in DG2.

	DG2	Miss A	Miss B	Miss C	Miss D	Miss E
Classical	0.875	0.833	1.000	1.000	0.833	0.833
Opsahl	0.611	0.500	0.667	0.667	0.600	0.714
Exclusive	0.600	0.500	1.000	0.500	0.500	0.750

*Example 2.5*

Evaluations of  $C$ ,  $C^*$ , and  $C^\circ$  at DG2 (Table 1) are illustrative: Each pair of women differ by at least one statistic, implying that they all occupy structurally distinct neighborhoods; none of the statistics, however, distinguishes them all. While  $C^*$  and  $C^\circ$  take lower values than  $C$ , the rankings of the actors are loosely correlated. Of particular interest are Mss. B and C, whom  $C^*$  and  $C$  do not distinguish but who take opposite values of  $C^\circ$ . At Miss B, the 4-path  $(A, 3, B, 4, E)$  is an open wedge to  $C^*$  but not a wedge at all to  $C^\circ$ ; at Miss C, the 4-path  $(D, 1, C, 5, E)$  is as a wedge to both  $C^*$  and  $C^\circ$  but only closed to  $C^*$ .

$C^*$  attributes high TC to Miss C because her friends remain better-connected when she is removed from the network, *while the events she attended remain*. In contrast,  $C^\circ$  attributes high TC to Miss B because her friends remain better-connected when she is removed from the network *along with the events she attended*. The statistic  $C^*$  thus detects TC that relies in part on inclusive events, which  $C^\circ$  does not.

The remainder of this section comes with a warning that the labeling schemes for triad nodes vary by context: Canonical triads  $\text{Tr}_{\mu w}$  have actors  $p, q, r$  such that  $w_{pq} \geq w_{qr} \geq w_{pr}$  (and unlabeled events);  $W$  and  $X$  adopt the schemes  $v_0, v_1, \dots$  for  $P_4$  and  $C_6$  from Sec. 1.2; and triads in larger ANs are scheduled at (ordered) triples of actors  $(i, j, k)$  with events  $a, b, \dots$

**Axioms** Sec. 1.2 delineated three goals for a new clustering coefficient: account for event size, as  $C^*$  does but  $C$  does not; further account for repeat group attendance, as neither  $C$  nor  $C^*$  do; and weight actors equally, as  $C$  does but  $C^*$  does not. This section wraps these desiderata into four axioms on  $\hat{C}$ . These are not suited to all settings, but they do help organize the myriad statistics that arise from Def. 2.4.

The first two axioms capture important features of  $C^*$ . In order to prevent single events from forming closed wedges,  $C^*$  is defined using only injections, from  $\mathcal{F}$ ; Axiom 1 requires that  $\mathcal{E}$  include induced injections (though not all injections). In order to allow distinct events to contribute distinct wedges,  $C^*$  removed the very strong congruence relation

<sup>10</sup> Some of these turn out to be the same statistic; for example, assuming  $X = \text{Tr}_{(1,1,1),0}$ , the category choices  $\overline{\mathcal{F}}/\simeq$  and  $\overline{\mathcal{F}}/\approx$  both yield  $C^\circ$ .

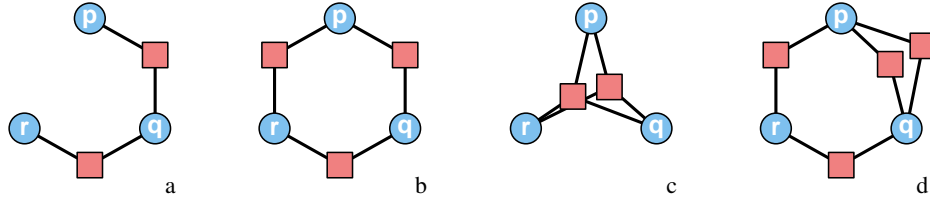


Fig. 5: Four AN triads from the axiomatic analysis: (a)  $\text{Tr}_{(1,1,0),0}$ , isomorphic to  $W$ ; (b)  $\text{Tr}_{(1,1,1),0}$ , isomorphic to  $X$ ; (c)  $\text{Tr}_{(0,0,0),2}$ , from the discussion of Axiom 4; and (d)  $\text{Tr}_{(2,1,1),0}$ , from the discussion of Lemma 2.14.

$\approx$  imposed on the morphisms of  $C$ ; Axiom 2 allows equivalences only when events are at least structurally equivalent.

*Axiom 1 (Induced injections)*

All induced injections (hence all isomorphisms) are morphisms (i.e.  $\mathcal{C}$  contains  $\overline{\mathcal{T}}$ ).

*Axiom 2 (Structural equivalence)*

The images of an (event) node under congruent morphisms are structurally equivalent (i.e.  $\sim$  is no stronger than  $\simeq$ ).

The last two axioms address the concerns raised with  $C^*$ . Axiom 3 addresses the problem of weighting by admitting at most one wedge at any ordered triple. Axiom 4 addresses the influence of bicliques by attacking their symptom: the counterintuitive way that each actor of a triad can have a wedge with none of the wedges being closed, which is not possible under  $C$ . The idea is that two wedges with different centers “hook together” (overlap) at their shared “side” (pair of actors), closing each other, which is here called “wedge buckling”. (Imagine rotating either “open” triad in Fig. 1 by  $120^\circ$  and overlaying it with itself.)  $C^*$  violates this idea, for example at Fig. 5c.<sup>11</sup>

*Axiom 3 (Equal representation)*

At each ordered triple there exists exactly one of the following: no wedge, one open wedge, or one closed wedge.

*Axiom 4 (Wedge buckle)*

If wedges exist at two ordered triples with different centers in a triad, then both are closed.

**Theorems** Three useful properties follow from certain subsets of the axioms: two triadic formulations of  $\hat{C}$ , which aide conceptualization and computation (Thms. 2.6 and 2.7), and one characterization (Thm. 2.8). The proofs constitute the next section.

*Theorem 2.6 (Census formulation)*

For each triad  $\text{Tr}_{\mu w}$ , write the numbers

$$\begin{aligned} F_{\mu w} &= |\text{Hom}_{\mathcal{C}/\sim}(W, \text{Tr}_{\mu w}) \setminus \text{Hom}_{\mathcal{C}/\sim}^X(W, \text{Tr}_{\mu w})| \\ S_{\mu w} &= |\text{Hom}_{\mathcal{C}/\sim}^X(W, \text{Tr}_{\mu w})| \end{aligned}$$

<sup>11</sup> One could instead simply impose as an axiom the restriction of wedges and alcoves to exclusive events; Axiom 4 provides an alternative framing for the problem.

of open and closed wedges, respectively, at  $\text{Tr}_{\mu w}$ . If  $\widehat{C}$  is defined using Eq. 4, then

$$\widehat{C}(G) = \frac{\sum_{\mu, w} s_{\mu w}(G) S_{\mu w}}{\sum_{\mu, w} s_{\mu w}(G) (F_{\mu w} + S_{\mu w})}. \quad (6)$$

Thm. 2.6 decomposes the rate-of-closure calculation into a ratio of motifs, according to the distribution of triads in  $G$ . The theorem proves useful in implementing the various global statistics, which may then be computed via arithmetic on the full census.

*Theorem 2.7 (Binning formulation)*

Assume Axioms 1, 3, and 4. Then the triads of  $G$  can be binned into subsets  $S_\emptyset(G)$ ,  $S_W(G)$ , and  $S_X(G)$  according as they contain none, two open, or six closed wedges; and

$$\widehat{C}(G) = \frac{3|S_X(G)|}{|S_W(G)| + 3|S_X(G)|}. \quad (7)$$

Thm. 2.7 generalizes the simple triad census description of  $C$  in Sec. 1.  $C^*$  does not satisfy these criteria, but  $C^\circ$  does: It is recoverable from the structural triad census as  $C^\circ = 3 \times (t_{30} + t_{31}) / (t_{20} + t_{21} + 3 \times (t_{30} + t_{31}))$ .

*Theorem 2.8 (Existence and uniqueness)*

There exist unique choices of  $X$ ,  $W$ ,  $\mathcal{C}$ , and  $\sim$  that satisfy Axioms 1, 2, and 3. Moreover, these choices also satisfy Axiom 4. Under them, Eqs. 4 and 5 both produce  $C^\circ$ .

Thm. 2.8 characterizes those specializations of  $\widehat{C}$  that satisfy every axiom.  $C^\circ$  turns out to be the unique such statistic; any alternative to  $C^\circ$  still expressible in terms of Def. 2.4 comes at the cost of at least one axiom. At the heart of Thm. 2.8 lies the tension between Axiom 2 and Axiom 3. The former forces different types of wedges to be treated differently, and the latter allows only one of these types to figure into the formula.

## 2.4 Proofs

**Triadic formulations** A different batch of lemmas leads up to each of the second two theorems, and Thm. 2.7 also depends on Thm. 2.6. To simplify the notation, in this section let  $\text{Hom}$  (with no subscript) denote the unspecified  $\text{Hom}_{\mathcal{C}/\sim}$ .

*Proof of Thm. 2.6*

The wedges  $\text{Hom}(W, G)$  can be partitioned according to which triad of  $G$  contains their images. The triads of  $G$  are, in turn, partitioned by the full census. Since the morphisms counts are fixed for isomorphic triads,

$$\widehat{C} = \frac{\sum_{H \subseteq G} |\text{Hom}^X(W, H)|}{\sum_{H \subseteq G} |\text{Hom}(W, H)|} = \frac{\sum_{\mu, w} \left( \sum_{\text{Tr}_{\mu w} \cong H \subseteq G} |\text{Hom}^X(W, \text{Tr}_{\mu w})| \right)}{\sum_{\mu, w} \left( \sum_{\text{Tr}_{\mu w} \cong H \subseteq G} |\text{Hom}(W, \text{Tr}_{\mu w})| \right)} = \frac{\sum_{\mu, w} s_{\mu w} \times S_{\mu w}}{\sum_{\mu, w} s_{\mu w} \times (F_{\mu w} + S_{\mu w})},$$

where  $H \subseteq G$  ranges over the triads of  $G$ .  $\square$

*Lemma 2.9*

Assume Axiom 1.

- i. If  $\text{Tr}_{\mu w}$  has an alcove, then every ordered triple of  $\text{Tr}_{\mu w}$  has an alcove.
- ii. Given actors  $i, j, k \in G$ , there is an openness-preserving bijection between the wedges of  $i, j, k$  and those of  $k, j, i$ .

Part i follows from the symmetry of  $X$ : Whatever the order of the actors, the structure of the triad is the same. Part ii follows analogously from the more limited symmetry of  $W$ , which allows  $v_0, v_1$  to be interchanged with  $v_4, v_3$  with no effect on the structure. (See Fig. 5a,b.)

*Proof*

For i, pick  $\psi \in \text{Hom}(X, \text{Tr}_{\mu w})$  and suppose  $\psi$  takes  $v_0, v_2, v_4$  to  $i, j, k$ . Pick any permutation  $\pi \in S_3$  so that  $\pi(i, j, k)$  is an arbitrary ordered triple in  $\text{Tr}_{\mu w}$ , and let  $\rho_\pi : X \rightarrow X$  be the isomorphism taking  $v_0, v_2, v_4$  to  $\pi(v_0, v_2, v_4)$ , which by Axiom 1 is a morphism. The composition  $\psi \circ \rho_\pi : X \rightarrow \text{Tr}_{\mu w}$  is then a morphism that takes  $v_0, v_2, v_4$  to  $\pi(i, j, k)$ .

For ii, let  $\rho : W \rightarrow W$  be the isomorphism on  $W$  that exchanges  $v_0$  and  $v_4$ , which is a morphism by Axiom 1. Composition with  $\rho$  assigns any wedge  $\phi : W \rightarrow G$  that sends  $v_0, v_2, v_4$  to  $i, j, k$  to a wedge  $\phi \circ \rho$  that sends  $v_0, v_2, v_4$  to  $k, j, i$ . Moreover, since  $\rho \circ \rho$  is the identity morphism on  $W$ , another composition with  $\rho$  takes  $\phi \circ \rho$  back to  $(\phi \circ \rho) \circ \rho = \phi \circ (\rho \circ \rho) = \phi$ . Composition with  $\rho$  thus pairs up the wedges of the triad  $i, j, k$  centered at  $j$  (no wedge is paired with itself). If such a wedge  $\phi$  factors through  $X$  as  $\phi = \psi \circ \iota$ , then  $\phi \circ \rho$  factors through  $X$  as  $\phi \circ \rho = (\psi \circ \iota) \circ \rho = \psi \circ (\iota \circ \rho) = \psi \circ (\rho' \circ \iota) = (\psi \circ \rho') \circ \iota$ , where  $\rho' : X \rightarrow X$  is the isomorphism on  $X$  that exchanges  $v_0$  and  $v_4$ . Thus  $\phi$  and  $\phi \circ \rho$  are open or closed together.  $\square$

The next two lemmas push the binning scheme of Thm. 2.6 from triads to ordered triples. The simplicity of Eq. 7 comes from the fixed number of possible wedges (one for each ordered triple; Axiom 3) and the symmetries between them (Lemma 2.9 and Axiom 4).

*Lemma 2.10*

Assume Axioms 1 and 4. Then, if a triad has two wedges with different centers, then every ordered triple in the triad has an alcove.

*Proof*

By Axiom 4, such a triad has a closed wedge, hence an alcove. By Lemma 2.9i, it then has an alcove at every ordered triple.  $\square$

*Lemma 2.11*

Assume Axioms 1, 3, and 4. Then each triad has exactly one of the following: no wedges, two open wedges, or six alcoves.

*Proof*

Each triad contains six ordered triples, which by Axiom 3 have at most one wedge each. Lemma 2.9ii requires that the wedges centered at any one actor either do not exist, are both open, or are both closed. Lemma 2.10 implies that, if two ordered triples with different centers have wedges, then all six have closed wedges. Thus the possible distributions of wedges among the six ordered triples are none, a pair of open wedges (at the same center), and six closed wedges.  $\square$

*Proof of Thm. 2.7*

Thm. 2.6 provides Eq. 6, which respects triad classes. Lemma 2.11 implies that either  $S_{\mu w} = F_{\mu w} = 0$ ,  $S_{\mu w} = 0$  and  $F_{\mu w} = 1$ , or  $S_{\mu w} = 3$  and  $F_{\mu w} = 0$  for every triad class. Binning these classes into  $S_\emptyset$ ,  $S_W$ , and  $S_X$ , respectively, achieves the result.  $\square$

**Characterization** The characterization theorem takes place over three steps: First, the three assumed axioms only allow wedges and alcoves with no inclusive events ( $\overline{\mathcal{T}}$ ). (This makes Axiom 4 unnecessary.) Second, the equal representation of Axiom 3 requires that any wedges at the same ordered triple of actors are congruent ( $\approx$ ), but when inclusive events are ignored the weaker relation  $\simeq$  is enough. This limits the options to the two formulations in Def. 2.4 under the category  $\mathcal{T} / \simeq$ . Third, these formulations agree under certain conditions, which turn out to be satisfied under  $\mathcal{T} / \simeq$ .

*Lemma 2.12*

Assume Axioms 1, 2, and 3. Then any wedge or alcove is an induced injection.

*Proof*

The only way for a wedge or alcove to not be an induced injection is for it to send some event to an inclusive event. Suppose the alcove  $\psi : X \rightarrow G$  sends  $v_0, v_1, v_2, v_3, v_4, v_5$  to  $i, d, j, e, k, f$ , where at least one of the events  $d, e, f$  is inclusive to the triad at  $i, j, k$ . ( $d, e$ , and  $f$  need not be distinct.) If  $d$  or  $e$  is inclusive, then  $\psi \circ \iota : W \rightarrow G$  is a wedge with an inclusive event. If only  $f$  is inclusive, then let  $\rho : X \rightarrow X$  be the isomorphism sending  $v_0, v_1, v_2, v_3, v_4, v_5$  to  $v_2, v_3, v_4, v_5, v_0, v_1$ , so that the composition  $\psi \circ \rho \circ \iota : W \rightarrow G$  sends  $v_0, v_1, v_2, v_3, v_4$  to  $j, e, k, f, i$ . By Axiom 1,  $\psi \circ \rho \circ \iota$  is a wedge with an inclusive event. It is enough, therefore, to prove the result for wedges.



Fig. 6: From the proof of Lemma 2.12: (a) the image of  $\psi : W \rightarrow G$  and (b) the necessary subgraph of  $G$  containing (a).

So suppose the wedge  $\phi : W \rightarrow G$  sends  $v_0, v_1, v_2, v_3, v_4$  to  $i, d, j, e, k$ , where at least one of  $d$  and  $e$  is inclusive to the triad at  $i, j, k$ . Obtain  $G'$  from  $G$  by adding events  $f$ , attended only by  $i$  and  $j$ , and  $g$ , attended only by  $j$  and  $k$ . (See Fig. 6.) The subgraph inclusion  $\sigma : G \rightarrow G'$  is an induced injection, hence by Axiom 1 a morphism. Then the composition  $\sigma \circ \phi : W \rightarrow G'$  is a wedge. The graph map  $\phi' : W \rightarrow G'$  sending  $v_0, v_1, v_2, v_3, v_4$  to  $i, f, j, g, k$  is an induced injection since  $f$  and  $g$  are exclusive events, so by Axiom 1  $\phi'$  is also a wedge—at the same ordered triple as  $\sigma \circ \phi$ . Axiom 2 implies that these wedges are incongruent, which contradicts Axiom 3. Thus  $\phi$  cannot exist.  $\square$

*Lemma 2.13*

Assume Axiom 3. Then  $\sim$  is at least as strong as  $\approx$  on the wedges and alcoves.

*Proof*

The claim is that any two wedges or alcoves on the same ordered triple of actors are congruent. If they were not, then Axiom 3 would be violated.  $\square$

The *pullback*  $\iota^* : \text{Hom}(X, G) \rightarrow \text{Hom}(W, G)$  sends any alcove  $\psi \in \text{Hom}(X, G)$  to the wedge  $\psi \circ \iota : W \rightarrow G$ . To understand Lemma 2.14, note that the image of  $\iota^*$  is in  $\text{Hom}^X(W, G)$ —that is, each such  $\psi \circ \iota$  factors through  $X$  (via the morphism  $\psi$  began with).

*Lemma 2.14*

Eqs. 4 and 5 yield the same statistic if and only if  $\iota^*$  is injective.

This lemma is not satisfied, for instance, by the category  $\widetilde{\mathcal{T}}/ =$  underlying  $C^*$ : The wedge  $\phi : W \rightarrow \text{Tr}_{(2,1,1),0}$  (Fig. 5d) sending  $v_0, v_2, v_4$  to  $v_2, v_4, v_0$  can be closed by either of the events shared by  $v_0$  and  $v_2$ .  $C^*$ , defined using Eq. 4, counts this as one closed wedge. Its counterpart  $\widehat{C}$ , defined using Eq. 5, however, counts two alcoves, one for each choice of event—that is,  $\phi$  factors through  $X$  in two ways. (Under this statistic, in fact,  $\widehat{C}(\text{Tr}_{(2,1,1),0}) = \frac{6}{5}$ .)

*Proof*

Given  $\phi \in \text{Hom}^X(W, G)$ , by definition there exists  $\psi \in \text{Hom}(X, G)$  such that  $\phi = \psi \circ \iota$ ; thus, in any case,  $\iota^*$  has image  $\text{Hom}^X(W, G)$ . The second condition therefore amounts to  $\iota^*$  being a bijective correspondence between its domain  $\text{Hom}(X, G)$  and its range  $\text{Hom}^X(W, G)$ . Since  $\iota^*$  is surjective and its domain and range are finite, this is true if and only if the domain and range have equal size. Since the denominators of Eqs. 4 and 5 are equal, this is true if and only if the formulations are equal, unless both are undefined. This occurs only when  $\text{Hom}(W, G)$  is empty, in which case both  $\text{Hom}(X, G)$  and  $\text{Hom}^X(W, G)$  are also empty.  $\square$

*Proof of Thm. 2.8*

Lemma 2.12 implies that wedges and alcoves are induced injections. By Axiom 1, all of these are morphisms. As far as Def. 2.4 is concerned, then,  $\mathcal{C}$  is  $\overline{\mathcal{T}}$ .

Lemma 2.13 implies that the congruence relation  $\sim$  is no weaker than  $\approx$ . Since the events of two wedges or alcoves at the same ordered triple must be exclusive, hence structurally equivalent in the triad, the relations  $\simeq$  and  $\approx$  have the same effect in this case;  $\mathcal{C}/\sim$  is  $\overline{\mathcal{T}}/\simeq$ . This establishes uniqueness.

For the auxiliary claim, suppose  $\psi, \psi' \in \text{Hom}_{\overline{\mathcal{T}}/\simeq}(X, G)$  are incongruent. By the choice of  $\overline{\mathcal{T}}$ , their respective images of  $v_1, v_3, v_5$  must be exclusive. If  $\psi, \psi'$  agree on all three actors, then, by the choice of  $\simeq$ , they are congruent. So  $\psi, \psi'$  must disagree on some actor; say  $\psi(v_0) \neq \psi'(v_0)$ . This implies that  $\psi \circ \iota(v_0) = \psi(v_0) \neq \psi'(v_0) = \psi' \circ \iota(v_0)$ , hence that  $\iota^*(\psi) \neq \iota^*(\psi')$ . Thus,  $\iota^*$  is injective. By Lemma 2.14, both formulations of Def. 2.4 produce the same statistic.

It remains to verify that  $C^\circ$  actually satisfies each axiom; this is left to the reader.  $\square$

### 3 Empirical analyses

This section applies triadic tools, including  $C$ ,  $C^*$ , and  $C^\circ$ , to three empirical networks. Sec. 3.1 assesses the clustering coefficients as measurement instruments, by comparing



their performances on the empirical networks. The assessments consider reliability, validity, redundancy, and practicality, and are illustrated in two case studies. Sec. 3.2 performs triadic analyses of the empirical networks, using the census and the clustering coefficients. The analyses draw upon and extend concepts from previous studies (see Sec. 1.1), including strong triadic closure, brokerage, and influence.

### 3.1 Instrumentation

**Data** The analyses employ three empirical networks: The social activity attendance network DG1 comes from another table in the same study as above (Davis *et al.*, 1941), and has seen extensive use as a test case for node classification and community detection techniques (Freeman, 2003). A subset of interlocking directorates data, from a study of corporate philanthropy in Minneapolis–St. Paul (Galaskiewicz, 1985; Wasserman & Faust, 1994), constitute GWF. Finally, MR refers to the collaboration network constructed from the *Mathematical Reviews* bibliographic database, which is maintained by the American Mathematical Society, over the years 1985–2008. These networks are constructed from a range of types and volumes of social interaction data and have appeared in previous studies that provide checks and comparisons for the present work. Two (DG1 and MR) have time-labeled events.<sup>12</sup>

Table 2: Structural censuses of DG1, GWF, and two intervals of MR. The column indicates the presence (1) or absence (0) of an inclusive event; the row indicates the number of pairs of actors who attend at least one exclusive event.

	DG1		GWF		MR (1985-7)		MR (2005-7)	
	0	1	0	1	0	1	0	1
0	0	17	0	284	80,747,526,018,836	17,275	725,892,036,097,769	76,558
1	39	240	266	886	4,721,138,210	8,611	38,496,757,064	51,599
2	146	253	452	521	133,630	2,014	909,505	15,185
3	45	76	130	61	886	129	5,585	1,055

Table 2 presents the structural censuses of the networks. The higher-order structure lost in projection lives mostly in the second column of each census. Several differences between DG1 and GWF, on one hand, and MR, on the other, are apparent: MR is far larger, with triads concentrated among the less-connected; “symmetric exclusive” triads ( $t_{30}$ , see Ex. 2.3) make up a minuscule fraction, undercut only by that of “symmetric complete” triads ( $t_{31}$ ). In contrast, DG1 and GWF have remarkably similar profiles: the event-free triads number  $t_{00} = 0$ , and the largest tallies occupy a northeast–southwest diagonal band away from the least and most connected types. This indicates that the smaller networks are more uniformly connected, with fewer poorly-connected actors. This difference likely

<sup>12</sup> DG1 is assumed to consist of events spanning nine months (Freeman, 2003); however, whereas the study took place over two years, other orderings are not impossible.

reflects non-uniformity in the coverage of researchers in MR (Lee & Cunningham, 2014), e.g. as equally prolific researchers on the periphery of mathematics appear less frequently in MR (Brunson *et al.*, 2014).

The editors assign to each publication one primary and any number of secondary Mathematical Subject Classification (MSC) codes from a hierarchical scheme. At the coarsest level, publications are binned into 64 groups (for instance, algebraic geometry, partial differential equations, and astronomy and astrophysics). For the assessments, 64 subnetworks are constructed by partitioning the literature by primary classification. Of these, 39 satisfy the following inclusion criteria over each adjacent 3-year interval from 1985–7 to 2006–8: the literature is not empty; each of  $C$ ,  $C^*$ ,  $C^\circ$ , and  $D$  is defined; and no two of these statistics are simultaneously zero. Since their curation and construction are systematic, differences in structure among these networks should only reflect differences in the cultures of research publication and limitations of MR coverage. (Nonetheless, size and density are known to influence measures of TC.)

**Criteria** While the statistics surveyed in Sec. 1.2 are hopefully intuitive, it is not yet clear that they are useful instruments.<sup>13</sup> This section assesses the local and global definitions of  $C$ ,  $C^*$ , and  $C^\circ$  on the basis of stability, concurrent validity, discriminability (meant to reflect practicality), and distinguishability (non-redundancy). The assessments are performed on three samples: the 18 actors of DG1, the 26 actors of GWF, and the 39 disciplines of MR (along adjacent 3-year intervals). The criteria are conceptualized and assessed as follows:

- An instrument is *stable* if it yields similar measurements of the same subject at different times. Stability is assessed, on pairs of values at the same MR discipline at adjacent intervals, as the proportion  $\frac{SSM}{SST}$  of the variation in the values accounted for by the pairing in a one-way analysis of variance (Altman & Bland, 1983).
- Both  $C^*$  and  $C^\circ$  are hypothesized to measure properties of graphs that can also be measured in other ways: As mentioned in Sec. 1.2, an alternative correction to  $C$  for event size in ANs is the quotient of  $C$  by its expected value  $C_{\text{rand}}$  on an equivalent random bipartite graph.<sup>14</sup> Sec. 1.2 also suggested that  $C^\circ$  may measure dynamic TC, defined as  $D$ . The *concurrent validity* of each measure shall be assessed as its coefficient of determination  $R^2$  with its alternative (Kimberlin & Winterstein, 2008).
- Two instruments designed to measure distinct properties shall be called *distinguishable* if they yield divergent values on the same subjects. Whereas the coefficient of determination between these values gives their concurrent validity, the remaining proportion of variance,  $1 - R^2$ , shall assess their distinguishability.
- An instrument is *discriminable* if its values in practice are dispersed throughout its theoretical range (Comin *et al.*, 2015). (Sec. 1 criticized  $C$  for having low discriminability on ANs.) Discriminability is assessed as the variance  $s^2$  of an instrument's

<sup>13</sup> Strictly speaking, the “instrument” that assigns a clustering coefficient to a social network includes the collection of sociometric data and the construction of the bipartite graph as well as the graph-theoretic calculation and the device that performs it; only the calculation is meant here.

<sup>14</sup> Here  $C_{\text{rand}}$  is calculated two ways: For the smaller networks DG1 and GWF, take the mean (local) values of  $C$  across 1000 randomly generated bipartite graphs having the same actor and event degree sequences (Chen *et al.*, 2005; Admiraal & Handcock, 2008). For the MR subnetworks, use the asymptotic approximation (Newman *et al.*, 2001).

values for a sample of subjects; the standardized values  $4s^2$  are reported, so that discriminability theoretically ranges from 0 (all values equal; statistic is useless) to 1 (values evenly split between 0 and 1; statistic perfectly dichotomizes the subjects). A statistic whose values follow a Gaussian distribution centered at 0.5 with standard deviation 0.25 (and cut off at the 95% thresholds) has discriminability just under  $\frac{1}{4}$ , while one whose values are uniformly distributed has discriminability  $\frac{2}{3}$ .

On MR, each assessment is performed on the pooled values across all intervals. For instance, each statistic's stability is computed on  $39 \times 7 = 273$  ordered pairs of values.

Table 3: Evaluations of three clustering coefficients taken over actors (DG1 and GWF) or subnetworks (adjacent 3-year intervals of MR).

	Classical			Opsahl			Exclusive		
	DG1	GWF	MR	DG1	GWF	MR	DG1	GWF	MR
Stability			0.781			0.403			0.457
Validity				0.622	0.296	0.113	0.058		0.399
Dist. (Classical)				0.950	0.940	0.999	0.492	0.732	0.924
Dist. (Opsahl)							0.915	0.592	0.948
Discriminability	0.005	0.013	0.047	0.051	0.050	0.026	0.205	0.224	0.001

**Results** The test results constitute Table 3. (Non-meaningful or redundant cells are left empty. Plots for each assessment are included in the supplement.)  $C$  is by far the most stable of the statistics ( $\frac{SSM}{SST} = 0.78$ ), with less than half of the variation in  $C^*$  and  $C^\circ$  each interval accounted for by the previous. Tests of validity were inconsistent.  $C^*$  was highly correlated with  $C/C_{\text{rand}}$  across the women of DG1, but much less so across the CEOs of GWF and the disciplines of MR. Conversely,  $C^\circ$  accounted for 40% of the variance in  $D$  across the disciplines but none across the women. Some heteroskedasticity is also visible in the plots of  $C^\circ$ . There is strong evidence here that these instruments are closely related, but only in certain limited settings.

The three statistics are highly distinguishable; at worst,  $C$  explains half of the variance in  $C^\circ$  across the women of DG1 ( $1 - R^2 = 0.49$ ). This, residual plots reveal, is due to a consistent negative relationship.  $C$  and  $C^*$  are poor discriminants, but on the actors of the smaller networks  $C^\circ$  takes values nearly as distributed over  $[0, 1]$  as the hypothetical cut-off Gaussian. This makes sense in light of the higher average rates of TC in DG1 and GWF; by comparison, the many highly-connected triads of MR are overwhelmed by the more partially-connected, which  $C$  is better-equipped to discriminate among (and does). Overall, the assessments lend some legitimacy to the uses of  $C$ ,  $C^*$ , and  $C^\circ$  in the next section, but more persuasive assessments of single-value network statistics would be helpful.

*Example 3.1*

Consider the TC of the women who constitute DG1 (Table 4, with structural equivalents Olivia and Flora represented by Olivia. Centrality scores will be used in Sec. 3.2. The

Table 4: Measures of local triadic closure and centrality in DG1.

	Classical	Opsahl	Exclusive	Dynamic	TwoWalk	Eigenvector	TwoWalkCorrected
Evelyn	0.897	0.767	0.448	0.576	0.319	0.335	0.015
Laura	0.962	0.842	0.487	0.692	0.286	0.309	0.023
Theresa	0.897	0.752	0.145	0.650	0.358	0.371	0.013
Brenda	0.962	0.839	0.450	0.692	0.292	0.313	0.021
Charlotte	1.000	1.000	1.000	1.000	0.154	0.168	0.014
Frances	0.962	0.869	0.778	0.000	0.198	0.209	0.011
Eleanor	0.962	0.796	0.531	0.692	0.220	0.228	0.008
Pearl	0.933	0.646	0.467	0.636	0.187	0.180	-0.007
Ruth	0.897	0.670	0.328	0.650	0.242	0.236	-0.006
Verne	0.897	0.674	0.393	0.576	0.231	0.218	-0.013
Myra	0.933	0.714	0.556	0.273	0.204	0.187	-0.017
Katherine	0.933	0.770	0.536	0.273	0.237	0.220	-0.017
Sylvia	0.897	0.746	0.300	0.576	0.292	0.277	-0.015
Nora	0.897	0.838	0.663	0.725	0.281	0.264	-0.017
Helen	0.897	0.816	0.661	0.611	0.215	0.201	-0.014
Dorothy	0.933	0.541	0.467	0.000	0.143	0.131	-0.012
Olivia	1.000	0.581	1.000	1.000	0.088	0.070	-0.019

supplement contains the table for GWF). Partitioning and core–periphery algorithms tend to identify Pearl, Ruth, and Verne as intergroup bridges or peripheral group members (Freeman, 2003), though in terms of classical TC their neighborhoods are unremarkable. In contrast, these women exhibit the lowest Opsahl TC of the group, and two (Ruth and Verne) are among the three with lowest exclusive TC. These observations attest to the greater discriminability of these statistics.

Pearl, however, has exclusive TC on par with several women in the cores of the two communities (Evelyn, Laura, and Dorothy). Theresa and Sylvia, on the other hand—who are usually placed near the cores of their respective groups within DG1, rather than toward the periphery with Ruth and Verne—show lower exclusive TC. This is due to the high number of events (8 and 7) these women attended. It may be that the study window omitted events attended by their neighbors in their absence, though both women attended events as early as March and as late as September, making this less likely; or it may be that these women played distinctive networking roles in their respective groups, to which traditional algorithms are not sensitive (see Sec. 3.2).

### Example 3.2

A previous study of MR (Brunson *et al.*, 2014) compared two subnetworks, constructed via a nearly even partition of primary MSCs into “pure” and “applied”.<sup>15</sup> The analysis of TC used  $C$  and  $C/C_{\text{rand}}$ ; the time series are reproduced in Fig. 7 (“Classical” and “BipartiteCorrected”). While  $C$  revealed persistent properties of MR, e.g. that the applied research

<sup>15</sup> The partition is coarse and provisional, but reveals a real difference between the research cultures; these subnetworks displayed consistently and characteristically different behavior.

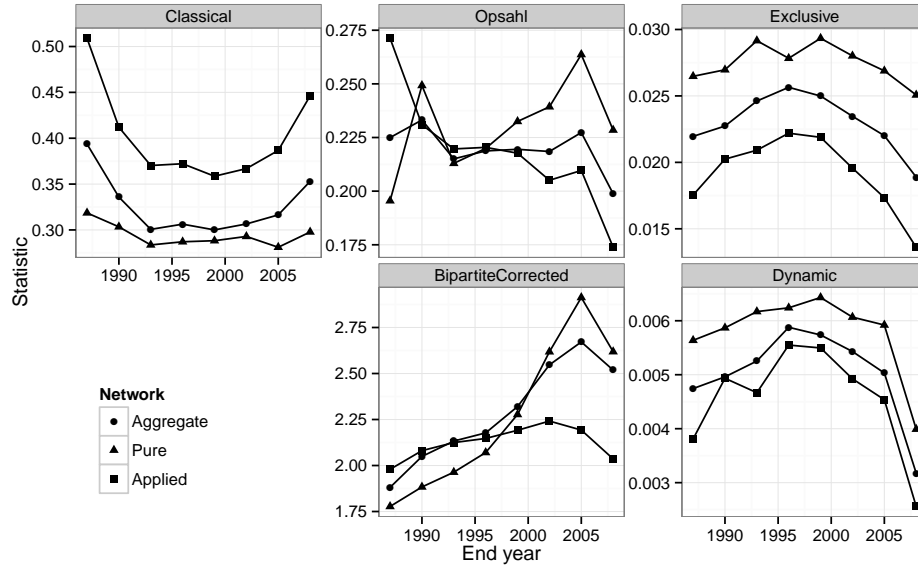


Fig. 7: Three global clustering coefficients and alternative measures for two, on the aggregate, pure, and applied MR networks along adjacent 3-year intervals.

community saw more classical TC than the pure,  $C/C_{rand}$  revealed discordant trends in pure and applied research. Both statistics arguably discriminated well, and certainly they were distinguishable from each other.

Fig. 7 also includes time series for  $C^*$  and  $C^\circ$ . The three trajectories of  $C^*$  mimic those of  $C/C_{rand}$  up to a linear transformation; the rates of change are clearly least in the pure network and greatest in the applied. More impressive is the stark resemblance between  $C^\circ$  and  $D$ , up to scale.  $C^*$  and  $C^\circ$  both are less discriminating than  $C$  in absolute terms, though all three are clearly distinguishable. Like  $C$ ,  $C^\circ$  measures a persistent difference between the research cultures: Pure research is better-characterized by exclusive (or dynamic) TC than applied. The negative relationship between  $C$  and  $C^\circ$  is evident here: the relative values of  $C^\circ$  are inverted from those of  $C$ , both in the ordering of the networks and in the concavity of the trends.

### 3.2 Triadic closure in affiliation networks

**Strong triadic closure** In social networks with ties of different strengths, the STC hypothesis predicts that, when two pairs of actors in a triad are *strongly* tied, then the third pair will tend to be at least *weakly* tied (Granovetter, 1973). Investigators have formalized and tested this principle in a variety of ways, often in terms of the frequency, duration, or intimacy of relations, or of the proportion of relations above some threshold of strength (Freeman, 1992). One conversion approach to STC in ANs is therefore to apply these methods to a weighted projection.

The full triad census offers a direct approach: Within a triad, it makes sense to infer stronger ties between actors from exclusive events than from inclusive events, consistent

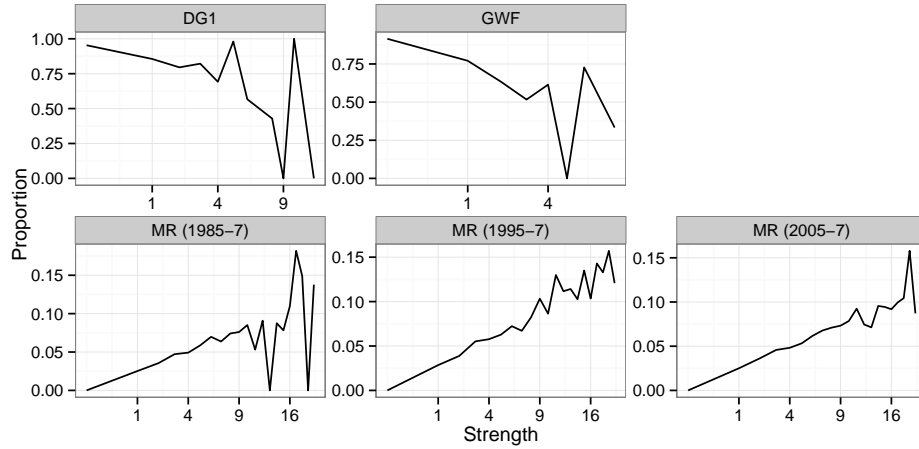


Fig. 8: Conditional probability  $\Pr(\mu_3 + w > 0 \mid \mu_1 \times \mu_2 = s)$  of a weak tie versus tie strength  $s$ , up to  $s = 20$ . (Note the square root-scale horizontal axis.)

with the principle that higher-attendance events foster weaker pairwise connections (Gupte & Eliassi-Rad, 2012). Accordingly, take the *wedge strength* of the ordered triple  $(i, j, k)$  to be the number of 4-paths along exclusive events from  $i$  through  $j$  to  $k$ , and take  $i$  and  $k$  to be (at least) *weakly tied* if there is any 2-path between them. Thus, the triple  $(p, q, r)$  in the triad  $\text{Tr}_{\mu_w}$  have wedge strength  $\mu_1 \times \mu_2$  and are weakly tied if  $\mu_3 + w > 0$ . STC shall be measured in an AN as the probability of a weak tie conditional on wedge strength.<sup>16</sup>

Fig. 8 presents the conditional probabilities for DG1, GWF, and MR over three evenly-spaced 3-year intervals, using a square-root scale on the horizontal axis. In DG1 and GWF, increasing wedge strength is associated (albeit noisily) with a lower rate of weak tie formation, in defiance of STC. In contrast, STC in MR is well-modeled by the proportionality

$$\Pr(\mu_3 + w > 0 \mid \mu_1 \times \mu_2 = s) \propto s^{\frac{1}{2}}. \quad (8)$$

Furthermore, though STC makes no predictions about the proportion of ties between actors who have no neighbors in common (the case  $s = 0$ ), in MR this case is accurately extrapolated from the pattern across wedges of positive strength.

**Connectedness and constraint** The STC hypothesis is intimately tied to the study of brokerage, in that connections among an actor  $i$ 's neighbors can be thought to constrain  $i$ 's potential to broker between them (Burt, 1992). Constraint is formulated as a product of  $i$ 's investment in connecting with their neighbors and the connectedness of these neighbors with each other. The local clustering coefficient provides a simple model of constraint: If  $i$  has  $d$  neighbors, each  $j$  of whom is tied to  $d(j)$  of  $i$ 's other neighbors, then the constraint

<sup>16</sup> An alternative measure is the expected number of events attended by  $i$  and  $k$ , conditioned on the wedge strength of  $(i, j, k)$ . The results in MR, not reported, are similar to those shown.

on  $i$  due to  $j$  can be defined as

$$c(i, j) = \frac{1}{d} \times \frac{d(j)}{d-1} = \frac{d(j)}{d(d-1)},$$

with total constraint  $c(i) = \sum_j c(i, j) = C(i)$ . The equivalent formulation

$$c(i, j) = \frac{|\{\text{wedges at } i \text{ w/ } j\}|}{|\{\text{wedges at } i\}|} \times \frac{|\{\text{closed wedges at } i \text{ w/ } j\}|}{|\{\text{wedges at } i \text{ w/ } j\}|} = \frac{|\{\text{closed wedges at } i \text{ w/ } j\}|}{|\{\text{wedges at } i\}|} \quad (9)$$

generalizes neatly to the terms of Def. 2.4. Thus the family of local clustering coefficients may be viewed as a family of alternative measures of constraint in ANs.<sup>17</sup>

This presents an opportunity to explore the relationship between connectedness and constraint. As originally defined, constraint decreases with neighborhood size, holding network density constant. A subtle change in definition, from a focus on proportional investment to one on marginal investment, instead produces polynomial *growth* in constraint due to a strong interaction effect with local density. In both theoretical (Szabó *et al.*, 2003) and empirical (Ravasz *et al.*, 2002; Vázquez, 2003) studies, the classical clustering coefficient exhibits the power law relationship

$$C_\ell \propto \ell^{-1}. \quad (10)$$

This may appear to conform to the former definition of constraint, but it actually concerns variation in local density. The family of measures encoded in Def. 2.4 may likewise be expected to behave differently, depending on the variety of TC they measure.

Taking  $\mathcal{C}$  to be  $\mathcal{T}$  and taking the trivial quotient by  $=$  effectively weights the local connectivity of  $i$ , as measured by the wedge count at  $i$ , by the number of  $i$ 's neighbors and the multiplicity of  $i$ 's shared events with them, moderated by the extent of overlap of these events among the neighbors. As a measure of constraint, then,  $C_\ell^*$  is highly sensitive to compounding constraint by multiple events, even between the same small subset of  $i$ 's neighbors. In contrast,  $C_\ell^\circ$  is sensitive only to pairs of  $i$ 's neighbors with at least one exclusive common event each (due to the restriction to  $\overline{\mathcal{T}}$ ) and is equally sensitive to constraints on  $i$ 's strategic position with respect to any such pair (Thm. 2.7). That is,  $C_\ell^*$  measures constraint weighted according to the strengths of the relationships (multiplicity of events) between  $i$  and two of their neighbors, while  $C_\ell^\circ$  measures constraint in the form of channels of exchange, hidden from  $i$ , between neighbors having their own exclusive channels with  $i$ .

Fig. 9 depicts  $C_\ell$ ,  $C_\ell^*$ , and  $C_\ell^\circ$  on MR, taken over the same three intervals as in Sec. 3.2.<sup>18</sup>  $C_\ell$  follows the expected power law-shaped curve, which persists over time. In contrast, the long-term trend in  $C_\ell^*$  is upward, and exhibits large fluctuations with persistent peaks (e.g. at  $\ell = 12$  and  $\ell = 24$ ), an expected artifact of biclique proliferation.<sup>19</sup>  $C_\ell^\circ$  mimics  $C_\ell$ : The long-term trend is downward and concave, and the fluctuations are modest and

<sup>17</sup> This should be compared cautiously to previous approaches that conditioned bipartite clustering coefficients on node degree (Lind *et al.*, 2005; Opsahl, 2013), rather than on a definition-specific wedge count.

<sup>18</sup> Scatterplots of values in DG1 and GWF are included in the supplement.

<sup>19</sup> Whenever  $n \geq 3$  and  $m \geq 2$ , the biclique  $K_{n,m}$  yields, for each of its actors  $j$ , pairs of neighbors and  $m(m-1)$  ordered pairs of events they share with  $j$ , resulting in  $(n-1)(n-2) \times m(m-1)$  4-paths

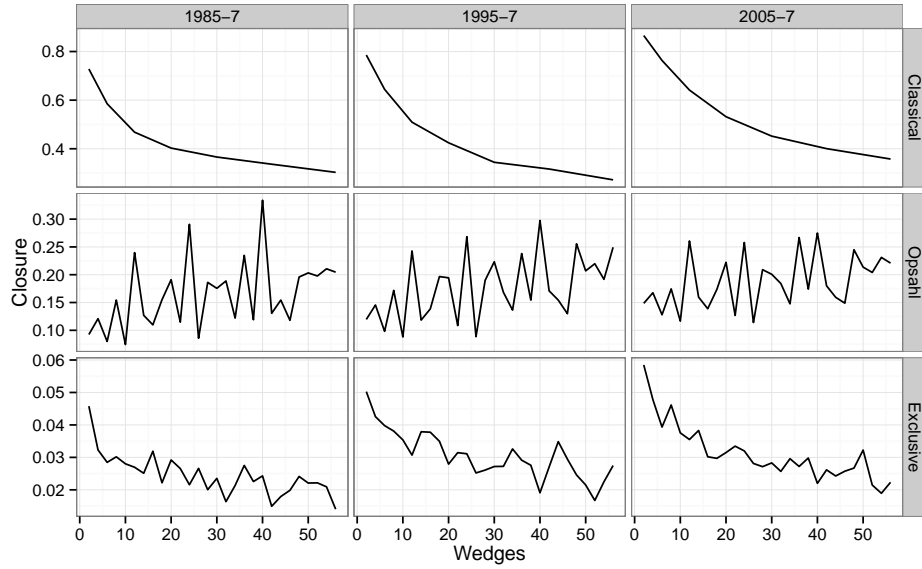


Fig. 9: Three wedge-dependent local clustering coefficients in MR. Note that  $C_\ell$  is only defined when  $\ell = k(k-1)$  for some integer  $k$ .

transient. Thus, in the world of research collaboration, the strengthening of one's (existing) collaborative ties may have a positive effect on the strengths of ties among one's collaborators; while the accumulation of new, mutually-exclusive ties is associated with fewer, on average, ties among them from which oneself is excluded.

Under the assumption that multiple shared events compound and interact to produce many multiple brokerage opportunities, the associated measure of constraint  $C_\ell^*$  may compound enough in kind to outpace it. In these terms, it is not necessarily to  $i$ 's advantage to accumulate neighbors through attendance at common events. In contrast, the constraint  $C_\ell^\circ$  imposed by exclusive channels among  $i$ 's neighbors diminishes with increased brokerage opportunities through  $i$ 's own exclusive channels. As in the classical case, therefore, it is unambiguously to  $i$ 's advantage to maintain many neighbors through mutually exclusive channels. These results demonstrate the range of possible behaviors for a custom measure of constraint, and the importance of specifying the brokerage patterns of interest.

**Constraint and influence** Like early conceptions of constraint, the preceding analysis focused on the structure of an actor  $i$ 's neighborhood. Yet much importance has also

centered at  $j$ . When  $m \geq 3$ , each of these is closed. Thus, any otherwise untied actor in a copy of  $K_{n,m}$  contributes the atypically high value  $C^*(j) = 1$  to the mean  $C_\ell^*$ , where  $\ell = (n-1)(n-2) \times m(m-1)$ . These values  $\ell = (3-1)(3-2) \times 3(3-1) = 12$ ,  $\ell = (3-1)(3-2) \times 4(4-1) = 24$ ,  $\ell = (4-1)(4-2) \times 3(3-1) = 36$ , and  $\ell = (3-1)(3-2) \times 5(5-1) = 40$  correspond to the highest peaks of  $C_\ell^*$  up to  $\ell = 56$ . Two clustering coefficients based on  $\widetilde{\mathcal{F}}/\simeq$  and  $\widetilde{\mathcal{F}}/\approx$  exhibited similarly expected fluctuations but decreased with wedge count. One based on  $\widetilde{\mathcal{F}}/=$  exhibited no such fluctuations and no long-term trend.



been placed on actors' positions within the entire network, as popular conceptions of centrality—closeness, betweenness, and eigenvector—attest. This last analysis attempts to discern whether the observed trade-offs are local or global phenomena, via a different extension of the same classical relationship.

Social influence is often measured by eigenvector centrality, based on the recursive principle that an actor accumulates influence through connections with other influential actors (Faust, 1997; Borgatti & Halgin, 2011). The eigenvector centrality of  $i$  can be expressed as the cumulative influence of  $i$  through walks (paths that may repeat nodes and edges) of at most some specified length; 1-walk centrality, for instance, equals node degree. This calculation can be inverted to produce a measure of influence through walks of *at least* some length (Bonacich, 1991): If the  $\ell$ -walk centrality scores of the nodes of an AN  $G$  constitute the vector  $\mathbf{c}_\ell = (c_\ell(1), \dots, c_\ell(n))$ , and the eigenvector centrality scores comprise  $\mathbf{c}_\infty$ ,<sup>20</sup> then the  *$\ell$ -walk-corrected centrality scores*, which may be positive or negative, constitute  $\mathbf{c}_\infty - \mathbf{c}_\ell$ . The actors' 2-walk centrality scores provide a measure of the local component of their influence that is self-contained, i.e. that does not depend on the measure of constraint. Their 2-walk-corrected centrality scores measure the global component.

Each of  $i$ 's neighbors is accessible to  $i$  via some 2-walk, so that  $c_2(i)$  may depend largely on the number of  $i$ 's neighbors. As the previous analysis revealed, however, how these 2-walks are counted is also important. The 2-walks from  $i$  are most closely related to the wedges of  $C^*$ , so it is reasonable to expect only a weak relationship between  $c_2(i)$  and  $C^*(i)$ . In contrast,  $C$  and  $C^\circ$  are insensitive to redundant 2-walks (from  $i$  to some neighbor  $j$ ). In order to decompose the relationship between constraint and influence, each clustering coefficient is considered versus each component (2-walk and 2-walk-corrected) of influence.

Fig. 10 plots the relationships for the CEOs of GWF. (Those for the women of DG1, included in the supplement, are similar.) Those with  $C^*$  are indeed weak, as are those with  $\mathbf{c}_\infty - \mathbf{c}_2$ . The standout is  $C^\circ$  versus  $c_2$ , and this holds too in DG1: In these small networks, at least, exclusive TC is associated with discernibly lower local influence. Specifically, an increase in 2-walk centrality of 0.1 corresponds to a decrease of 0.46 (GWF) or 0.31 (DG1) in  $C^\circ$ . The lack of any discernible relationship with 2-walk-corrected centrality suggests that the configuration of  $i$ 's neighborhood is only weakly, if at all, related to  $i$ 's extended influence.

#### 4 Conclusion

This study pursued a measure of triadic closure for affiliation networks, modeled as bipartite graphs, that controls for the proliferation of bicliques. Bicliques arise from attendance at multiple events by subsets of actors, which is unlikely to reflect the popular understanding of triadic closure. The need for such a measure follows from the sensitivity of existing measures to such structures, even those that control for the sizes of events. In addition to the proposed exclusive clustering coefficient  $C^\circ$ , the paper presented a classification scheme

<sup>20</sup> Here each  $\mathbf{c} = \mathbf{c}_\ell, \mathbf{c}_\infty$  is normalized so that  $\sum_i \mathbf{c}(i)^2 = 1$ .

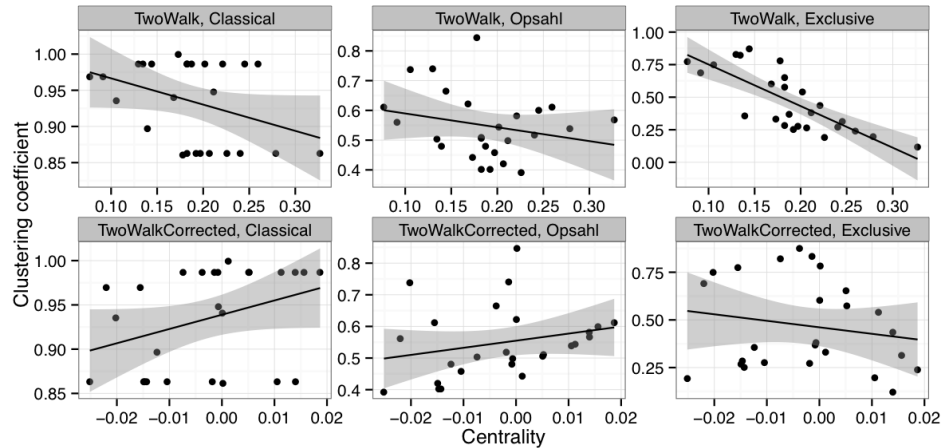


Fig. 10: Scatterplots of Opsahl and exclusive clustering coefficients versus 2-walk and 4-walk–corrected eigenvector centrality scores across actors in GWF. Least-squares regression lines and 95% confidence bands are overlaid.

for affiliation network triads and an axiomatic framework for defining affiliation network clustering coefficients.

An instrumental analysis found  $C^\circ$  to measure distinct properties from the classical  $C$  and the recent proposal  $C^*$ , and suggested that, in some settings,  $C^\circ$  approximates triadic closure as it is characterized over time. An investigation of several empirical affiliation networks revealed patterns of triadic closure much richer than could be inferred from the classical triad census and  $C$  applied to their actor projections. In the author’s judgment,  $C^\circ$  comes across as a useful counterpoint to  $C$ ; the two could be viewed as limiting cases between which other clustering coefficients like  $C^*$  interpolate (Saramäki *et al.*, 2007).

The study has several limitations, most notably the limited number of empirical (and lack of simulated) affiliation networks investigated, and the fact that these networks were constructed using different data collection methods. This leaves the conclusions drawn here open to challenge. Also, no fast algorithms were provided, and the implementations used were not designed for efficiency; applications of the tools described here to large networks will require both.

The tools suggest several other avenues for future work. The classification of affiliation network triads provides the basis for a state transition analysis, which may aide models of network evolution. Affiliation networks also exist with weighted edges, and the generic clustering coefficient described could be adapted, like its predecessor  $C^*$ , to this setting.

In summary, it is hoped that the present paper provides a useful framework for the triadic analysis of affiliation networks.

**Acknowledgments** The author thanks Ritchie C. Vaughan for suggesting this line of inquiry; to Miranda Lynch, Tina Eliassi-Rad, Roldan Pozo, Paola Vera-Licona, Linton Freeman, Kathy O’Hara, and Reinhard Laubenbacher for helpful conversations; to Barry Brunson and Pansy Waycaster for several rounds of proofreading; to four anonymous reviewers

for highly incisive and supportive comments; and to the Virginia Bioinformatics Institute, the AMS, and UConn Health for data and resources. This project builds upon work done by participants in the Summer 2010 REU in Modeling and Simulation in Systems Biology.

### References

- Admiraal, R., & Handcock, M. S. (2008). networksis: A package to simulate bipartite graphs with fixed marginals through sequential importance sampling. *Journal of statistical software*, **24**(8), 1–21.
- Altman, D. G., & Bland, J. M. (1983). Measurement in Medicine: The Analysis of Method Comparison Studies. *The statistician*, **32**(3), 307–317.
- Bonacich, P. (1991). Simultaneous group and individual centralities. *Social networks*, **13**(2), 155–168.
- Bondy, J. A., & Murty, U. S. R. (1976). *Graph theory*. New York: Elsevier Science Publishing.
- Borgatti, S. P., & Everett, M. G. (1997). Network analysis of 2-mode data. *Social networks*, **19**, 243–269.
- Borgatti, S. P., & Halgin, D. S. (2011). Analyzing affiliation networks. Pages 417–433 of: Scott, John, & Carrington, Peter J. (eds), *The Sage handbook of social network analysis*. London: Sage Publications Ltd.
- Breiger, R. L. (1974). The Duality of Persons and Groups. *Social forces*, **53**(2), 181–190.
- Brunson, J. C., Fassino, S., McInnes, A., Narayan, M., Richardson, B., Franck, C., Ion, P., & Laubenbacher, R. C. (2014). Evolutionary events in a mathematical sciences research collaboration network. *Scientometrics*, **99**(3), 973–998.
- Burt, R. S. (1992). *Structural holes: The social structure of competition*. Cambridge, MA: Harvard University Press.
- Carrino, C. N. (2006). *A study of repeat collaboration in social affiliation networks*. Ph.D. thesis, University Park, PA, USA. AAI3343661.
- Chen, Y., Diaconis, P., Holmes, S. P., & Liu, J. S. (2005). Sequential Monte Carlo Methods for Statistical Analysis of Tables. *Journal of the american statistical association*, **100**(469), 109–120.
- Comin, C. H., Silva, F. N., & da F. Costa, L. (2015). A framework for evaluating complex networks measurements. *EPL (Europhysics letters)*, **110**(6), 68002.
- Csardi, G., & Nepusz, T. (2006). The igraph software package for complex network research. *Interjournal, Complex Systems*, 1695.
- Davis, A., Gardner, B. B., & Gardner, M. R. (1941). *Deep south; a social anthropological study of caste and class*. Chicago: The University of Chicago Press.
- Davis, J. A. (1967). Clustering and structural balance in graphs. *Human relations*, **20**(2), 181–187.
- de Sola Pool, I., & Kochen, M. (1978). Contacts and influence. *Socnet*, **1**(1), 5–51.
- Easley, D., & Kleinberg, J. (2010). *Networks, crowds, and markets: Reasoning about a highly connected world*. New York, NY, USA: Cambridge University Press.
- Faust, K. (1997). Centrality in affiliation networks. *Social networks*, **19**(2), 157–191.
- Freeman, L. C. (1992). The sociological concept of "group": An empirical test of two models. *The american journal of sociology*, **98**(1), 152–166.
- Freeman, L. C. (2003). Finding social groups: A meta-analysis of the southern women data. Pages 39–97 of: Breiger, Ronald, Carley, Kathleen, & Pattison, Philippa (eds), *Dynamic social network modeling and analysis: Workshop summary and papers*. National Academies Press.
- Galaskiewicz, J. (1985). *Social organization of an urban grants economy: A study of business philanthropy and nonprofit organizations*. Orlando, FL: Academic Press.

- Glänzel, W., & Schubert, A. (2004). *Analyzing scientific networks through co-authorship*. Open Access publications from Katholieke Universiteit Leuven. Katholieke Universiteit Leuven.
- Granovetter, M. S. (1973). The strength of weak ties. *The american journal of sociology*, **78**(6), 1360–1380.
- Gupte, M., & Eliassi-Rad, T. (2012). Measuring tie strength in implicit social networks. *Pages 109–118 of: Proceedings of the 4th annual ACM Web Science Conference*. WebSci '12. New York, NY, USA: ACM.
- Harary, F., & Kimmel, H. J. (1979). Matrix measures for transitivity and balance. *J. Math. Sociol.*, **6**(2), 199–210.
- Hell, P. (1979). An introduction to the category of graphs. *Pages 120–136 of: Topics in graph theory (New York, 1977)*. Ann. New York Acad. Sci., vol. 328. New York Acad. Sci., New York.
- Holland, P. W., & Leinhardt, S. (1971). Transitivity in structural models of small groups. *Small group research*, **2**(2), 107–124.
- Kimberlin, C. L., & Winterstein, A. G. (2008). Validity and reliability of measurement instruments used in research. *Am j health syst pharm*, **65**(23), 2276–2284.
- Kreher, D. L., & Stinson, D. R. (1999). Combinatorial algorithms: Generation, enumeration, and search. *SIGACT news*, **30**(1), 33–35.
- Lee, C., & Cunningham, P. (2014). Community detection: effective evaluation on large social networks. *Journal of complex networks*, **2**(1), 19–37.
- Liebig, J., & Rao, A. (2014). Identifying influential nodes in bipartite networks using the clustering coefficient. *Pages 323–330 of: Proceedings of the tenth international conference on signal-image technology and internet-based systems*.
- Lind, P. G., González, M. C., & Herrmann, H. J. (2005). Cycles and clustering in bipartite networks. *Phys. rev. E*, **72**(Nov), 056127.
- Martin, T., Ball, B., Karrer, B., & Newman, M. E. J. (2013). Coauthorship and citation patterns in the Physical Review. *Phys. rev. E*, **88**(Jul), 012814.
- Mitchell, B. (1965). *Theory of categories*. Pure and Applied Mathematics, vol. 17. New York and London: Academic Press.
- Newman, M. E. J. (2001). Scientific collaboration networks. I. Network construction and fundamental results. *Phys. rev. E*, **64**(Jun), 016131.
- Newman, M. E. J. (2003). The structure and function of complex networks. *SIAM rev.*, **45**(2), 167–256 (electronic).
- Newman, M. E. J., Strogatz, S. H., & Watts, D. J. (2001). Random graphs with arbitrary degree distributions and their applications. *Phys. rev. E*, **64**(Jul), 026118.
- Opsahl, T. (2013). Triadic closure in two-mode networks: Redefining the global and local clustering coefficients. *Social networks*, **35**(2), 159 – 167. Special Issue on Advances in Two-mode Social Networks.
- R Development Core Team. (2008). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.
- Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N., & Barabási, A.-L. (2002). Hierarchical organization of modularity in metabolic networks. *Science*, **297**(5586), 1551.
- Saramäki, J., Kivela, M., Onnela, J.-P., Kaski, K., & Kertész, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. *Phys. rev. E*, **75**(Feb), 027105.
- Stanley, R. P. (2002). *Enumerative combinatorics*. Cambridge studies in advanced mathematics, no. v. 1. Cambridge University Press.
- Szabó, G., Alava, M., & Kertész, J. (2003). Structural transitions in scale-free networks. *Phys rev e*, **67**(5), 056102.
- Uzzi, B., & Spiro, J. (2005). Collaboration and creativity: The small world problem. *American journal of sociology*, **111**(2), 447–504.

- Vázquez, A. (2003). Growing network with local rules: Preferential attachment, clustering hierarchy, and degree correlations. *Phys. rev. E*, **67**(May), 056104.
- Wasserman, S., & Faust, K. (1994). *Social network analysis: Methods and applications*. Vol. 8. Cambridge University Press.
- Watts, D. J., & Strogatz, S. H. (1998). Collective dynamics of “small-world” networks. *Nature*, **393**(6684), 440–442.
- Wickham, H. (2009). *ggplot2: Elegant graphics for data analysis*. Springer New York.

### A Supplement

Figs. A 1–A 4 elaborate upon the scores in Table 3. Table A 1 is the counterpart, for GWF, to Table 4 in the main text. Fig. A 5 is the counterpart, for DG1 and GWF, to Fig. 9 in the main text, except that the ordered pair for every actor is plotted, rather than their wedge-dependent averages. Fig. A 6 is the counterpart, for DG1, to Fig. 10 in the main text.

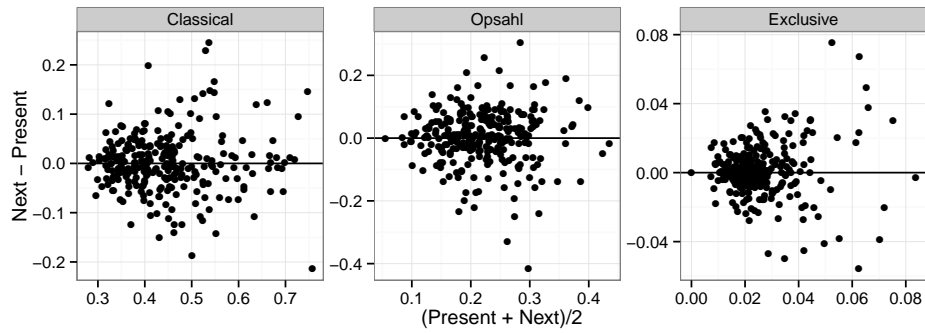


Fig. A 1: Mean-difference plots for values of  $C$ ,  $C^*$ , and  $C^o$ , taken across 39 subnetworks of MR over 7 pairs of adjacent intervals.

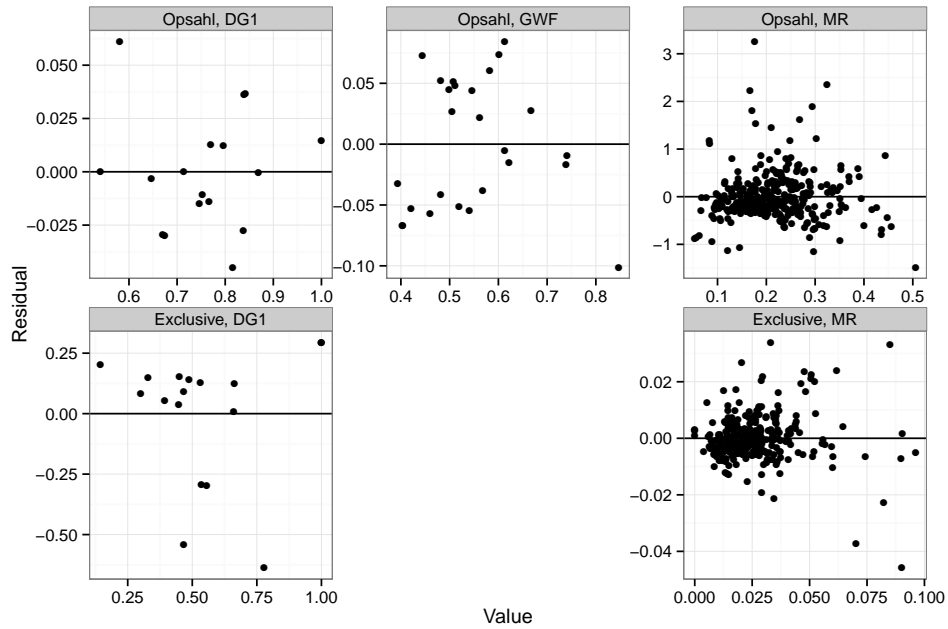


Fig. A 2: Residual plots for  $C/C_{\text{rand}}$  regressed on  $C^*$  and  $D$  regressed on  $C^o$ , taken across the women of DG1, the CEOs of GWF, and 39 subnetworks of MR over 8 intervals.

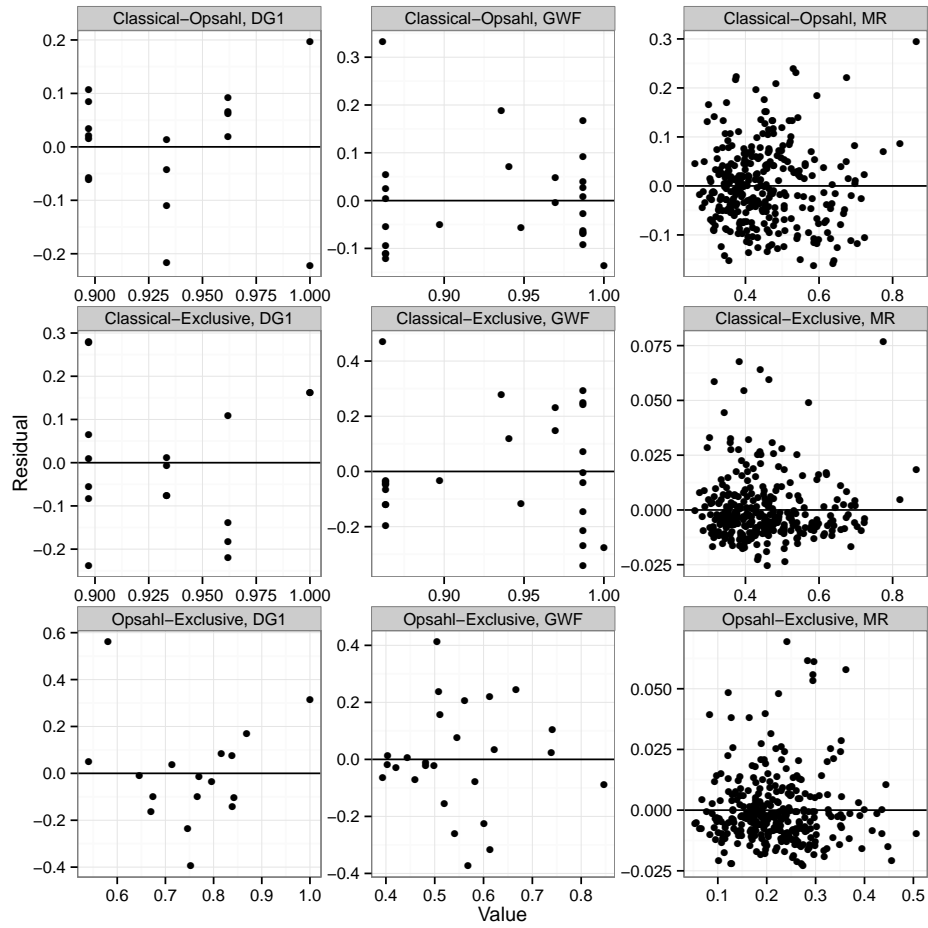


Fig. A 3: Residual plots for  $C^*$  regressed on  $C$ ,  $C^\circ$  regressed on  $C$ , and  $C^\circ$  regressed on  $C^*$ , taken across the women of DG1, the CEOs of GWF, and 39 subnetworks of MR over 8 intervals.

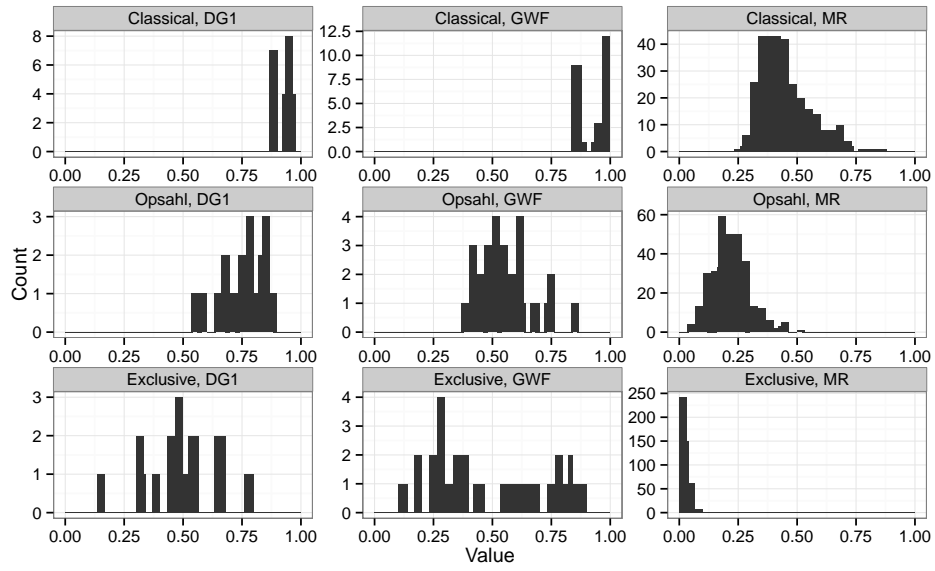


Fig. A 4: Histograms of values of  $C$ ,  $C^*$ , and  $C^\circ$ , taken across the women of DG1, the CEOs of GWF, and 39 subnetworks of MR over 8 intervals.

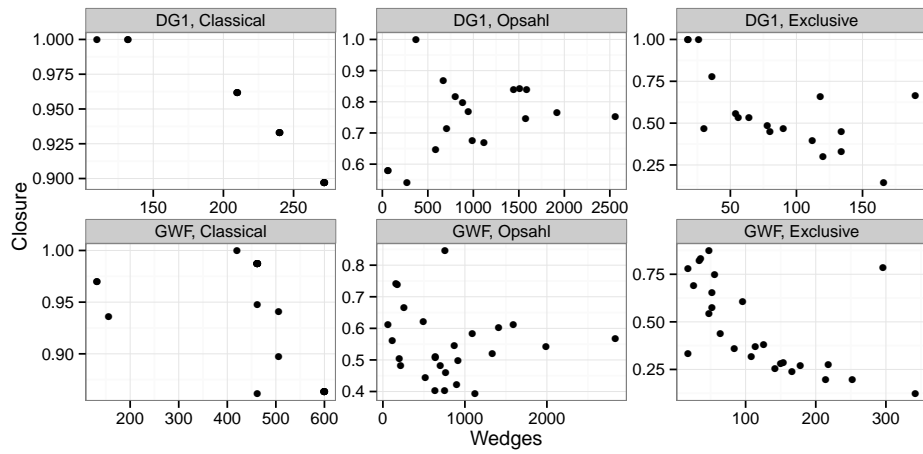


Fig. A 5: Four wedge-dependent local clustering coefficients in DG1 and GWF.



Table A 1: Measures of local triadic closure and centrality in GWF.

	Classical	Opsahl	Exclusive	TwoWalk	Eigenvector	TwoWalkCorrected
CEO1	0.863	0.403	0.254	0.192	0.178	-0.014
CEO2	0.897	0.481	0.357	0.139	0.127	-0.012
CEO3	0.987	0.741	0.833	0.130	0.128	-0.001
CEO4	0.987	0.546	0.542	0.202	0.213	0.011
CEO5	0.987	0.667	0.875	0.144	0.140	-0.004
CEO6	1.000	0.444	0.333	0.173	0.174	0.001
CEO7	0.863	0.460	0.280	0.197	0.187	-0.010
CEO8	0.970	0.561	0.692	0.091	0.069	-0.022
CEO9	0.936	0.739	0.750	0.106	0.086	-0.020
CEO10	0.987	0.505	0.824	0.135	0.127	-0.007
CEO11	0.987	0.481	0.368	0.188	0.187	-0.001
CEO12	0.970	0.613	0.778	0.077	0.061	-0.016
CEO13	0.863	0.421	0.270	0.207	0.192	-0.015
CEO14	0.863	0.568	0.123	0.327	0.341	0.014
CEO15	0.987	0.601	0.315	0.245	0.261	0.016
CEO16	0.948	0.499	0.381	0.212	0.211	-0.001
CEO17	0.987	0.613	0.241	0.260	0.278	0.019
CEO18	0.861	0.847	0.784	0.178	0.178	0.000
CEO19	0.863	0.393	0.196	0.226	0.201	-0.025
CEO20	0.863	0.541	0.198	0.279	0.289	0.011
CEO21	0.863	0.404	0.286	0.183	0.168	-0.015
CEO22	0.941	0.622	0.604	0.168	0.168	0.000
CEO23	0.987	0.582	0.438	0.221	0.235	0.014
CEO24	0.863	0.519	0.275	0.240	0.239	-0.002
CEO25	0.987	0.508	0.654	0.183	0.188	0.005
CEO26	0.987	0.511	0.577	0.183	0.188	0.005

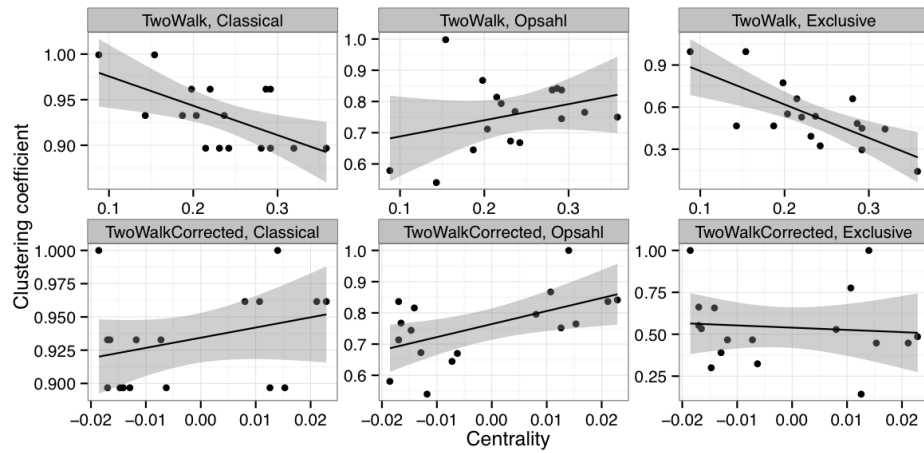


Fig. A6: Scatterplots of Opsahl and exclusive clustering coefficients versus 2-walk and 4-walk-corrected eigenvector centrality scores across actors in DG1. Least-squares regression lines and 95% confidence bands are overlaid.