

# Alliance Formation in a Multipolar World

Peter Devine<sup>y</sup>

George Washington University

Sumit Joshi<sup>z</sup>

George Washington University

Ahmed Saber Mahmud<sup>x</sup>

Virginia Polytechnic Institute and State University

May 2024

## Abstract

We propose a multilayer network approach to alliance formation. In a signed adjacency layer, agents are partitioned into clusters, with friendly relations within and hostile connections across clusters. Agents then form defensive collaborations in an alliance layer as follows: Agents in the same cluster form a nested split graph with degree inversely correlated to the level of hostility, and agents from disparate clusters with high-degree and low-hostility form cliques. Within cliques, agents from a cluster that is “intermediate” in terms of discord serve as a bridge to interconnect agents from more “extreme” clusters.

**Key words and phrases:** Alliance formation, signed graphs, nested split graphs, pairwise stability, cliques.

JEL Classification: C72, D74, D85

## 1 Introduction

This paper is a contribution to the literature on alliance formation under conflict. It explores the incentives of agents (individuals, groups or nations) to form defense alliances when they are embedded in a pre-existing network of bilateral adjacencies that are friendly, hostile or neutral, and for these adjacencies in turn to be revised following the formation of defense collaborations. This primitive non-empty network (equivalently *layer* or *graph*) of *adjacencies* is assumed to be a possible consequence of political, religious, ideological, cultural or historical factors.<sup>1</sup> It can be formally represented as a *signed* network in which a positive link between two

---

<sup>1</sup>We would like to express our sincere thanks to the editor, associate editor, and two anonymous referees whose comments

agents denotes friends, a negative link denotes enemies, and lack of a link denotes a neutral relationship.<sup>2</sup> While there is a large literature on network formation under conflict<sup>3</sup>, our point of departure is an explicit two-way interaction between the signed network of affinities and the network of defense *alliances* among friends to thwart potential conflicts with enemies. This interaction is examined through the lens of a *multilayer* network. Figure 1 depicts a multilayer network in which the base layer is the affinity network (denoted by  $H$ ) and the accompanying layer is the network of defense alliances (denoted by  $G$ ). The affinity layer is a signed network in which a solid line connecting two agents denotes friendship (a positive relationship), a dashed line denotes hostility (a negative relationship), and lack of a connection denotes a neutral (zero) relationship. The alliance layer is an unsigned network in which a (solid) line connecting two agents denotes a defense collaboration and the absence of a line implies no such collaboration.



Figure 1: A Multilayer Network

Our paper is motivated by the fact that the complex web of interlocking defense alliances that characterize the world today can best be understood as a multilayer network building up from base affinities. The period of the Cold War was characterized by an affinity network in which countries were broadly divided into an Eastern and a Western bloc based on opposing political ideologies. The corresponding alliance network was *bipolar*: the Eastern bloc formed the Warsaw Pact while the Western bloc formed NATO, with no overlap between the two security pacts. The fall of the Berlin Wall altered the affinity network with former Eastern bloc countries recalibrating their relationships with the Western bloc. The resulting alliance network was *unipolar* with former Warsaw Pact members such as Poland, Hungary, Bulgaria,

<sup>2</sup>Signed networks are discussed in Cartwright and Harary (1956), Davis (1967), and Easley and Kleinberg (2010, Chapter 5).

<sup>3</sup>Please see Bloch (2012) and Goyal et al. (2016) for an excellent description of the main lines of research on alliance formation under conflict.

Romania and the Czech Republic joining NATO. The current alliance network is sometimes described as *multipolar*, which is inaccurate since nations cannot be divided into mutually exclusive coalitions that jointly coordinate their actions as a set. Instead, we see nations forming alliances across continents. For example, the Economist<sup>4</sup> has noted that the United States has established bilateral alliances with Australia, Japan, Philippines, South Korea and Thailand in a hub-and-spoke network and quoted the prime minister of Japan, Kishida Fumio, as saying that promoting alliances among the spokes “will lead to the establishment of a *multilayered network*”

*Second*, under what circumstances will agents have an incentive to revise their relationships in the ac nity network? Specifically, what are the incentives of agents to mend fences with enemies and transform hostile relationships into friendships? *Third*, how will any change in the ac nity network impact defensive collaborations in the alliance network? *Fourth*, and finally, who are the agents that serve as “bridges” in the ac nity network to connect agents who would otherwise remain disconnected due to their mutual hostility?

We begin with a description of the architecture of the ac nity network  $H$  in the initial position. We assume that the distribution of positive and negative links is such that we can partition agents into non-empty *clusters* such that relationships within a cluster are friendly or neutral while relations across clusters are hostile or neutral. In the terminology of signed graphs following Davis (1967), a network with this particular distribution of positive and negative links is called *weakly balanced* (henceforth, simply *balanced*). In balanced ac nity networks, each cluster is composed of agents who are friends, friends of friends, friends of friends of friends... etcetera. Any links connecting agents across distinct clusters are always negative indicating that the agents are enemies. We assume that the partition of agents into clusters is a consequence of their disagreement over some norm. Agents within the same cluster subscribe to a common norm or core belief (for example, ideology, religion or politics) and thus any links that exist within the cluster are always friendly. Agents in distinct clusters differ in their perception of the norm and this dissonance implies that any links in  $H$  connecting an agent-pair from two separate clusters is always hostile. The norm or belief is captured by a scalar and thus permits classifying clusters as “close” or “distant” depending on the difference between their adoptive norms or beliefs.

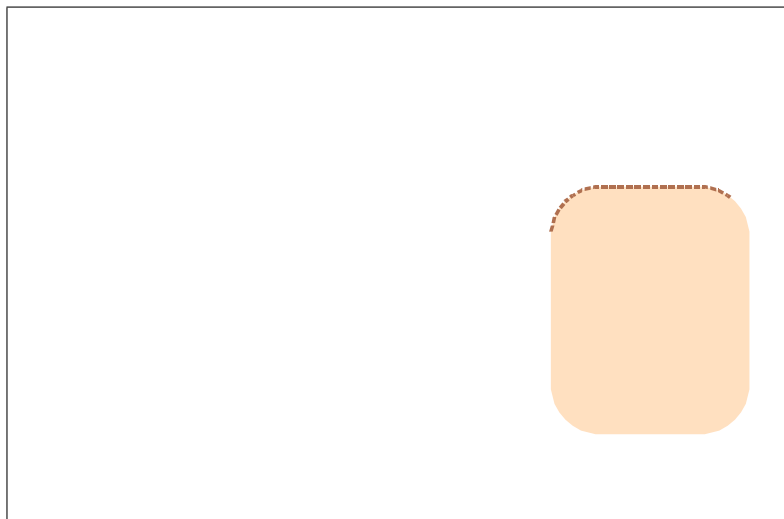


Figure 2: A Balanced Ac nity Network

Figure 2 illustrates a balanced ac nity network with three clusters. Each agent  $i$  in  $H$  is indexed by a *friendship* measure,  $f_i$ , which is the number of  $i$ 's friends minus the number of  $i$ 's enemies. The *higher* the value of  $f_i$ , the more friends agent  $i$  has relative to enemies, and thus the *lower* is the level of hostility



the most hostility in  $H$  is the least connected in  $G$  with a neighborhood contained in the neighborhoods of all other agents.

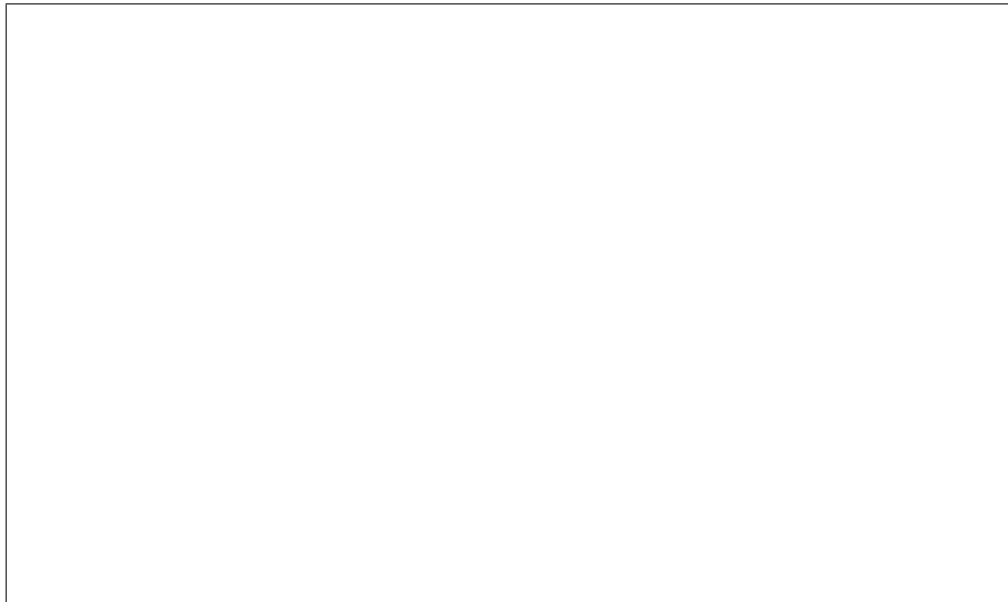


Figure 4: Architecture of the Alliance Network

Once the alliance network is formed, we allow agents to revisit their relationships in the affinity network. Therefore, we permit a two-way interaction between the affinity and alliance networks. Once agents have formed sufficient alliances within their own cluster, then the ensuing gains from these links can provide an incentive for well-connected agents in two separate clusters to change an existing hostile relationship or a neutral one in  $H$  into a friendly one. This will be particularly true if the difference between their perceived norms is sufficiently small. Of course, agents could also transform a neutral relationship into a hostile one but in our model there is no incentive to do so. Once these changes in the affinity network are implemented, then this revised (and potentially unbalanced) affinity network will spur a new round of alliances in the network  $G$ . Since new pathways of positive (friendly) links have been created across disparate pairs of clusters, two sufficiently well-connected agents from disparate clusters have an incentive to ally with each other. In particular, we show that *if two clusters are sufficiently close in their perceived norms, then sufficiently well-connected agents in the two clusters form a clique in the alliance network. A clique in  $G$  is a set of agents such that every pair of agents in the sets are mutually linked.* Thus, despite their dissonance over the norm, erstwhile hostile agents will have an incentive to ally if their disagreement over the norm is small.









Agents are assumed to belong to different clusters because of discord over what they believe should be the norm. The norm is captured by a scalar taking values over an interval  $[\underline{c}, \bar{c}]$ , where  $0 < \underline{c} < \bar{c} < 1$ . Agents *within* a cluster  $C \in \{C_0, C_1\}$ , irrespective of whether they are neutral or friends in  $H_0$ , subscribe to a common norm  $c_i \in [\underline{c}, \bar{c}]$ , and  $c_i \in C_0$  if  $c_i \in C_0$ . This norm is assumed immutable and does not change. The greater the difference,  $|c_i - c_j|$ , the more agents in clusters  $C_0$  and  $C_1$  differ in terms of core beliefs. We will define for  $i \in C_0$  and  $j \in C_1$ :

$$d_{ij}(H_0) = \frac{1}{1 + |c_i - c_j|} \quad (1)$$

We will use  $d_{ij}(H_0)$  as a *measure of discord* between clusters and suppress reference to  $H_0$  for brevity. If agents  $i$  and  $j$  belong to the same cluster, then  $d_{ij} = 1$  and there is no discord; if they belong to different clusters, then  $d_{ij} < 1$ . Thus, the greater the dispersion in subscribed norms, the lower the value of  $d_{ij}$ . Note that the measure of discord is a property of two clusters and not specifically of agents; in other words, for distinct agents  $i, j, k, l$  where  $i, k \in C_0$  and  $j, l \in C_1$ , we have  $d_{ij} = d_{kl}$ . Also note that  $d_{ij} = d_{ji}$ . It is important to note once again that the discord between agents is fixed with respect to their position in  $H_0$ . Even if subsequently two agents  $i$  and  $j$  from different clusters establish a friendly relationship, their mutual discord  $d_{ij}$  is not equal to 1, i.e., they are still not in consonance with respect to their respective subscribed norms.

Agents will be permitted to make limited changes to the primitive  $H_0$ . A pair of agents  $i$  and  $j$  can change the relationship from neutral or enemy to friend, by each side incurring a cost that captures the effort required to build the necessary trust. Thus, the formation of a friendly link requires *bilateral* consent of the pair of agents involved. The individual cost to agents  $i$  and  $j$  of converting  $h_{ij} \in \{0, 1\}$  to  $h_{ij}$

and let  $d_i^{(0)} = 0$  even if there are no isolated agents in  $G$ . The *degree partition* of  $G$  is denoted by  $D(G) = \{D_0(G); D_1(G); \dots; D_m(G)\}$ , where all agents in the element  $D_k(G)$  of the partition have the same degree  $d_k$ ,  $k \in \{0; 1; \dots; m\}$ . The definition of path and connectedness are defined analogous to the case of signed networks. A maximally connected subnetwork  $G^0$  in  $G$  is called a *component* of  $G$ . Given networks  $G$  and  $G^0$ , we will say that  $G$  is *denser* than  $G^0$  if  $G^0 \subset G$ . We will let  $G - ij$  (respectively,  $G + ij$ ) denote the network obtained from  $G$  by deleting (respectively, adding) the link  $ij$ .

An important network architecture that we will consider is a *nested split graph*. This network has the property that if  $d_i(G) < d_j(G)$ , then  $N_i(G) \subset N_j(G) \cup \{j\}$ . In other words, the neighborhood of a lower degree agent is contained within the neighborhood of a higher degree agent. Therefore, all allies of a less connected agent are also the allies of a more connected agent. Figure 3 illustrates this class of networks.

### 2.3 Gross Benefits from Alliances

Let  $Z$  denote the set of integers and consider the functions  $v: Z_+ \rightarrow R_+$  and  $w: Z_+ \rightarrow R_+$ . The function  $v$  captures return from own degree while  $w$  captures the return from the partner's degree. Suppose agent  $i$  with degree  $d_i$  forms a link with agent  $j$  with degree  $d_j$ . The incremental gross benefit to agent  $i$  from this link with agent  $j$  depends on the degree of both agents involved and is assumed to be given by:

the region dramatically reducing operating and maintenance costs for its own fleet. Furthermore, there are indirect benefits from aligning with a higher degree node in  $G$ . The AUKUS military alliance came about as Australia was set to join a looser and more transactional military industrial agreement with France. Australia’s post-cold war defense strategy concluded that “[Australia is] one of the most secure countries in the world...distant from the main centres of global military confrontation”<sup>10</sup>. Therefore a relatively inexpensive and limited agreement with France suited both. However, when Australia perceived China as a more present threat to Australia’s homeland, the country reneged on the agreement with France and opted for AUKUS with the two most central nodes in NATO and the western alliance at approximately five times the monetary cost and incurring significant obligations on its military autonomy and sovereignty. The benefits of joining the tripartite AUKUS with stronger ties between and emanating from each node were, *ceteris paribus*, significantly greater than a bilateral agreement with France.<sup>11</sup>

*Remark: (Separable versus non-separable gross benefits)* We have postulated an additively separable in degrees specification for gross benefits in (2). Such a separable specification allows a transparent exposition of the main results. A non-separable formulation, in which gross benefit to agent  $i$  from a link with agent  $j$  is more generally specified as  $b_{ij}$ , would also yield the same set of results under appropriate conditions on  $b_{ij}$ . We demonstrate this in Section 6.3.

## 2.4 Cost of Hostility

Recall that agents can only form an alliance in  $G$  with those who are friends or distant friends in  $H$  and that an alliance requires mutual consent. By forming an alliance, an agent incurs a cost of linking which is a function of the hostility faced by the potential partner in  $H$ . Letting  $c : Z \rightarrow \mathbb{R}_+$  denote this linking cost function, we will impose the following assumptions on  $c$ .

**Assumption (A.3):** For all  $z \in Z$ :

- (a)  $c(z + 1) < c(z)$ .
- (b)  $c(z) - c(z + 1) > c(z - 1) - c(z)$ .

Therefore, the cost to an agent is lower when the potential partner faces less hostility. Further, the cost reduction realized with a higher friendship partner is greater than with a lower friendship partner. With a link, each agent assumes some of the risks posed by the hostile relationships of the potential partner. These risks are consequently lower if each partner has more friends and less enemies. This also explains

<sup>10</sup>Protection by Projection, *The Economist*, April 25, 2023.

<sup>11</sup>Another example substantiating assumption (A.2) is the Nordic fighter fleet agreements signed in spring 2023. The Nordic nations (Norway, Sweden, Denmark, and Finland), agreed to pool resources creating an integrated air defense. The Scandinavian peninsula is a great example of this effect. Because the shortest distance from Russian air bases to the allied coast is from the north, each nation has a relatively small geographic slice of detection zones, but their entire geography is collectively exposed around the clock. Therefore, the increased number of participants in resource-sharing greatly expands the benefit to each individual member through burden sharing in time, distance and capacity by eliminating duplications in detection, early warning, rotating alert and surge units, and command and control infrastructure. A similar dynamic is also observed in the NATO VJTF, a rotating multinational task force kept on ready alert to spearhead an immediate counter attack to an invasion. All NATO members rotate units through the VJTF to defend the collective’s eastern perimeter.

why NATO does not permit admitting nations that are engaged in territorial disputes. A dispute indicates a high level of hostility towards that potential member and, therefore, a risk of conflict that commits the entire alliance. A more detailed look at this is NATO's admission of Finland, before the (presumptive) admission of Sweden. While both Finland and Sweden distrust Russia, only Sweden has some measure of hostile links with Turkey (also a NATO member). Finland, with its lower level of hostility, was therefore prioritized for admission to the military alliance (G) while the hostile link between Sweden and Turkey is being flipped in H.

Each agent will also face some cost due to its own hostile relations. We will let  $c_0 : Z \rightarrow \mathbb{R}_+$  denote the cost to an agent from these hostile relations and assume that  $c_0$  also satisfies conditions (a) and (b) of assumption A.3.

**Remark: (Costs that are functions of both hostility and degree)** The cost to agent  $i$  of linking with agent  $j$  is assumed to be  $c(i, j)$ . It can be argued that this cost to agent  $i$  could be lower if  $j$  had greater degree, i.e., the cost function should be  $c(i, j; d_j)$ , with  $c(i, j; d_j + 1) < c(i, j; d_j)$ . Our results would continue to obtain under this specification as well as demonstrated in Section 6.3. We note, however, that this effect is already captured in our basic model if we construe the net cost to agent  $i \in C(H_0)$  from linking with agent  $j$  as:

$$e_{i, j; j} = \begin{cases} c_0(j) & v_j; & j = i \\ c(i, j) & w_j; & j \in C(H_0) \\ c(i, j) & ijw_j & h_{ij}^+; & j \notin C(H_0) \end{cases} \quad (3)$$

Recalling assumption A.2(a),  $e_{i, j; j}$  is decreasing in  $j$ .

Finally, we impose a joint restriction on the cost of forming an alliance in G and the cost of transforming a relationship in H. This assumption bounds the reduction in alliance costs that can be achieved with agents changing their affinity relationship *within* a cluster from neutral to friendly. The logic is that agents are already friends or distant friends within a cluster. Thus, any reduction in alliance costs attained within a cluster by forming a more direct friendly relation is less than the cost of transforming an affinity relationship. This shifts the impetus of agents to revise links *outside* rather than inside the cluster in an affinity network.

**Assumption (A.4):** For all  $i \in Z$ :

$$c_0(i) - c_0(i + 1) <$$

Henceforth we will use the following notation to denote a unit increase in degree and friendship:

$$\begin{aligned} v(i) &= v(i + 1) - v(i) \\ c(i) &= c(i) - c(i + 1) \\ c_0(i) &= c_0(i) - c_0(i + 1) \end{aligned}$$

Note that we define  $c(i)$  and  $c_0(i)$  such that they are positive due to A.3(a).

## 2.5 Payoffs

In contrast to the contest function approach of the traditional literature, we adopt a reduced form additive specification of payoffs that reflect the tradeoffs present in the model. There are essentially four factors at play: (i) the “economies of scale” from allying with those who have high degree in  $G$ ; (ii) the cost of

cost of reaching out beyond their cluster. Specifically, if  $G$  (respectively,  $G^*$ ) is the stable alliance network under free riding (respectively, without free riding), then  $G^* \subseteq G$ .

### 3 Fixed Alliance Network

We begin our analysis with the case of a fixed alliance network  $H_0$ . We then examine the implications of this fixed alliance network on the topology of alliances in  $G$ . Therefore, we consider a one-way interaction between the alliance and alliance networks. This section can be construed as a *short run* analysis when the horizon is sufficiently small for agents to expect a change in relationships in the alliance network. We are assuming here that relationships (whether friend or enemy) embodied in the network  $H_0$  have taken time to coalesce. Within the time frame of the short run, new relations cannot be established in the alliance network. The incremental utility to agent  $i$  from forging an alliance in  $G$  with a member  $j$  in its own cluster is given by:

$$u_i(G + g_{ij}; H_0) - u_i(G; H_0) = [v_i(G) + 1 - v_i(G)] + w_j(G) + 1 - c_j(H_0)$$

We will use a definition of stability inspired by Jackson and Wolinsky (1996).

Definition (Pairwise-stability for monolayer networks): Given  $H_0$ , a network  $G$  is *pairwise-stable* if:

No agent  $i \in N$  has an incentive to unilaterally delete an existing link with agent  $j$  in  $G$ , i.e.,  $u_i(G - g_{ij}; H_0) - u_i(G; H_0) \leq 0$ .

No pair of agents  $i, j \in N$  who are unlinked in  $G$   $g_{ij} > 0$ .

g d()TJ/F68 10156t25H

### 3.1 The Basic Link Formation Game

We will examine link formation in  $G$  through a dynamic game inspired by Aumann and Myerson (1988). The advantage of this approach is that it selects one among potentially multiple pairwise stable networks.<sup>12</sup>



pair of networks in the set and no improving path leading to a network outside the set. We will show below (Theorem 1) that a closed cycle is not possible in our link formation game. Thus, the only outcome is convergence to a limit network  $G(H_0)$ .

**Theorem 1** *The basic link formation game converges to a limit network  $G(H_0)$  which is a pairwise-stable network.*

Therefore, Theorem 1 also shows the *existence* of a pairwise-stable network. The proof is based on the fact that no agent has an incentive to delete a link that it formed along an improving path in  $G$ . With deletions of links ruled out, cycles cannot emerge along an improving path. Therefore, since the number of network architectures are finite, the link formation game will converge to a pairwise-stable network.

**Remark: (Salient features of the basic game)** We note two facts about the dynamic game. *First*, we have the active agent deleting any unprofitable links and then proposing a new link to a passive agent. However, it is immaterial in our framework whether agents first delete links and then form a link, or first form a link and then delete links. This is because the incremental payoff from links that are formed will only increase by virtue of assumption A.2 as the degree of agents increase. Thus, as noted earlier, formed links are never subsequently deleted. *Second*, we allow the active agent to propose at most *one* link to a potential ally. We address in the next subsection the proposal of multiple links by an active agent.

### 3.2 The Pairwise-Stable Architecture of $G(H_0)$

Since friends and distant friends are contained within a cluster from assumption A.1, all alliances are between members of the same cluster. We will characterize the *intra*-cluster alliances formed by agents in  $G$  given  $H_0$ . Consider a given cluster  $C(H_0)$ , let  $I = \{i_1; i_2; \dots; i_n\}$  denote the set of agents arranged in increasing order of their index who belong to this cluster. Let  $|C(H_0)|$  denote the size of this cluster. For ease of exposition, let us assume without loss of generality that:

$$i_1 \quad i_2 \quad i_3 \quad \dots \quad i_n \tag{5}$$

with at 13336

The *friendship partition*,  $\mathcal{C}(H_0) = \{C_1(H_0); \dots; C_r(H_0); \dots; C_s(H_0)\}$ , is the collection of friendship classes in  $\mathcal{C}(H_0)$ .

Friendship classes will play an important role in our characterization result. All agents within the same friendship class face the same level of hostility. Agents belonging to a lower-index friendship class face greater hostility than agents belonging to higher-index friendship class. Thus, for example,  $i_j \in C_1(H_0)$  and  $i_n \in C_s(H_0)$ . Recalling the definition of a degree partition, let  $D(G) = \{D_0(G); D_1(G); \dots; D_m(G)\}$  denote the *degree partition* of agents belonging to  $\mathcal{C}(H_0)$  in the limit network  $G = G(H_0)$ . We will now examine how agents are distributed across this degree partition as a function of their friendship measure, i.e., their hostility level. Specifically, we will connect  $D(G)$  to the friendship partition  $\mathcal{C}(H_0)$ .

We will begin by elaborating on how link formation proceeds according to our dynamic game. Recall that (active) agents proceed in increasing order of their index and can only propose alliances with those who belong to their cluster. Therefore, we can consider how links are formed within any given cluster, say  $C_i(H_0)$ . The *first* active agent in  $C_i(H_0)$  to propose an alliance will be  $i_1$ . Note that at this stage, say  $G_0(\cdot)$ , no alliances have been formed and thus  $d_i(G_0(\cdot)) = 0$  for all  $i \in C_i(H_0)$ . Of course, if  $i_1$  is agent 1 in cluster  $C_1(H_0)$  who initiates the game, then  $G_0(\cdot) = G^e$ . The incremental payoff to agent  $i_1$  from proposing an alliance with agent  $i_k$  is  $(0) + w(1) - c_{i_1 i_k}$ . Since  $c_{i_1 i_k} = c_{i_1 i_j}$  for all  $i_k \in \{i_1, \dots, i_n\}$ , it follows that the most profitable alliance is with agent  $i_n$ . However, if this alliance yields a negative payoff to at least one agent, then the link will not be formed. We can of course have isolated



paribus choose the one with lower enmity. Thus, at each stage of the link formation game, preferential attachment implies that an agent facing lower hostility will have at least as many alliances as an agent facing greater hostility.

Another way to visualize the intra-cluster NSG architecture is as a “core-periphery” subnetwork composed of a hierarchical order of agents according to their degree. The peripheral agents are those who are connected only to the core agents but not among themselves; the core agents are connected to all other core agents and differ only with respect to the peripheral agents they are connected to. Recall the degree partition  $D(G) = \{D_0(G); D_1(G); \dots; D_m(G)\}$  within the cluster  $C(H_0)$  and let  $\lfloor x \rfloor$  denote the largest integer smaller than or equal to  $x$ . The peripheral agents arranged in increasing number of alliances and their set of allies are as follows:

Table 1: Peripheral Agents in  $C(H_0)$

Peripheral agents	Set of allies
$D_1(G)$	$D_m(G)$
$D_2(G)$	$D_m(G) [ D_{m-1}(G)$
$D_3(G)$	$D_m(G) [ D_{m-1}(G) [ D_{m-2}(G)$
$D_{\lfloor \frac{b}{2} \rfloor C}(G)$	$D_m(G) [ D_{m-1}(G) [ D_{m-2}(G) [ \dots [ D_{\lfloor \frac{b}{2} \rfloor C+1}(G)$

The smallest set of peripheral agents are those in  $D_1(G)$  connected only to the core agents in  $D_m(G)$  while the largest set of peripheral agents are those in  $D_{\lfloor \frac{b}{2} \rfloor C}(G)$  who are connected to all core agents. The core agents arranged in decreasing number of alliances are as follows:

Table 2: Core Agents in  $C(H_0)$

Core agents	Set of allies
$D_m(G)$	$D_1(G) [ D_2(G) [ D_3(G) [ D_4(G) [ \dots [ D_m(G)$
$D_{m-1}(G)$	$D_2(G) [ D_3(G) [ D_4(G) [ \dots [ D_m(G)$
$D_{m-2}(G)$	$D_3(G) [ D_4(G) [ \dots [ D_m(G)$
$D_{\lfloor \frac{b}{2} \rfloor C+1}(G)$	$D_{\lfloor \frac{b}{2} \rfloor C}(G) [ D_{\lfloor \frac{b}{2} \rfloor C+1}(G) [ \dots [ D_m(G)$

The agents in  $D_m(G)$  are core agents with the largest number of allies and they are connected to all agents – whether peripheral or core – in their cluster. Agents in  $D_{\lfloor \frac{b}{2} \rfloor C+1}(G)$  are core agents with the fewest number of allies and, while being connected to all core agents, are only allied with peripheral agents in the set  $D_{\lfloor \frac{b}{2} \rfloor C}(G)$ .

We now identify an interesting

## 4 Variable Affinity Network

So far we have kept the affinity network as fixed and examined its influence on the alliance network. However, it is possible that after harnessing sufficient economies of scale from their alliances, highly connected

## 4.1 The Augmented Link Formation Game

We now consider an augmented sequential link formation game that accommodates changes in both the acuity and alliance networks.

Given a non-empty primitive network  $H_0$ , link formation starts in the alliance network  $G$  starting from an empty network. The sequential process of link formation on this layer culminates in a limit network that we now denote as  $G_1 = G(H_0)$ .

The game now shifts to the network  $H$ . Players once again move sequentially in the order of their index starting from the state  $H_0^{(0)} = H_0; G(H_0); 1$ . The action set of the active agent in  $H$  is different from that in  $G$ . *First*, no links in  $H$  can be deleted. This is in accordance with our assumption that acuity relationships have matured bilaterally over a period of time and thus cannot be expunged unilaterally. *Second*, there is no incentive to convert a friend into an enemy because this makes an agent relatively unattractive as an ally to a potential partner. *Third*, by virtue of assumption (A.4), there is no incentive to revise a relationship within a cluster. An agent  $i$  can change an existing neutral relationship with agent  $j$  within the cluster to one of direct friends and the consequent increase in the friendship measure bestows a gain in own costs equal to  $c_0(i(H_0)) - c_0(i(H_0) + 1) < 0$ . This is consonant with our formulation that any transformation of acuity links is a precursor to forging alliances in the alliance network, and two agents within the same cluster do not have to resort to this intermediate step in order to connect in  $G$ . Therefore, the only choice we allow an active agent is to commit resources to convert a hostile or neutral relation *outside* the cluster into a friendly one.

Suppose agent  $i \in C_i(H_0)$  is the active agent. The active agent  $i$  can propose to an agent  $j \in C_o(H_0)$ ,  $i \neq j$ , with whom  $h_{ij} \in \{-1, 0, 1\}$  to change the relationship to a friend (i.e., to  $h_{ij} = +1$ ). Note that if another agent  $k \in C_i(H_0)$  had prior to  $i$ 's move already established a friendly relation with some agent, say  $l$ , in  $C_o(H_0)$ , then  $i$  has no incentive to make an overture to  $j \in C_o(H_0)$ . This is because a friendly path between clusters  $C_i(H_0)$  and  $C_o(H_0)$  has already been created in the acuity network through  $h_{kl} = +1$ . Thus, agent  $i$  can free ride on this link to form alliances in  $G$  with members of  $C_o(H_0)$  without having to first transform an acuity link with  $j \in C_o(H_0)$ .

Suppose, therefore, that when agent  $i \in C_i(H_0)$  is the active agent and proposes to agent  $j \in C_o(H_0)$ , then there is no friendly path connecting clusters  $C_i(H_0)$  and  $C_o(H_0)$ . This new relationship in the acuity network imposes a cost of  $c > 0$  but increases  $i$ 's friendship measure (lowers hostility level) to  $i(H_0) + 1$ . This increase in the friendship measure confers two benefits to agent  $i$ . *First*, by virtue of assumption A.3(a), it decreases  $i$ 's own costs:

$$c_0(i(H_0)) = c_0(i(H_0)) - c_0(i(H_0) + 1) > 0 \tag{9}$$

*Second*, by virtue of assumption A.3(b), it increases  $i$ 's own costs:

The active agent  $i$  chooses a potential partner  $j$  from another cluster with whom the sum of (9) and



Theorem 2 *The augmented link formation game converges to a limit  $(G = G(H); H = H(G))$  which is pairwise-stable.*

## 4.2 The Pairwise-Stable multilayer Network

We now characterize the pairwise-stable multilayer network  $(G; H)$ . Consider the augmented link transformation game when it moves from layer  $G_1$  to layer  $H_1$ . Consider any two clusters  $C = C(H_0)$  and  $C_0 = C_0(H_0)$ . Let  $i \in C = C(H_0)$  and  $j \in C_0 = C_0(H_0)$  be the most connected agents in their respective clusters (with the highest index agent chosen in case of a tie). Note from Proposition 2 that degree correlates positively with friendship, and thus these two agents are also the ones facing the lowest hostility in their respective clusters. Thus, as the following lemma indicates, these agents are the most likely candidates to transform their relationship to a friendly one since they realize the highest incremental utilities within their cluster from such a transformation in  $H$  and a subsequent alliance in  $G$ . Let  $D = D(G_1) = \{D_0(G_1); D_1(G_1); \dots; D_m(G_1)\}$  denote the degree partition of agents belonging to  $C = C(H_0)$  in the network  $G_1$  and define  $D^0(G_1)$  analogously. Also, let  $s = s(C = C(H_0))$  (respectively,  $s^0 = s(C_0 = C_0(H_0))$ ) denote the highest friendship class in  $C = C(H_0)$  (respectively,  $C_0 = C_0(H_0)$ ). From Proposition 2 we know that  $s = s(C = C(H_0)) = D_m(G_1)$  and  $s^0 = s(C_0 = C_0(H_0)) = D_m^0(G_1)$ .

Lemma 1 *Consider any two clusters  $C = C(H_0)$  and  $C_0 = C_0(H_0)$  and let  $i \in C = C(H_0) \setminus s = s(H_0)$  and  $j \in C_0 = C_0(H_0) \setminus s^0 = s^0(H_0)$ . For any  $k \in C = C(H_0)$  and  $l \in C_0 = C_0(H_0)$ :*

$$\begin{aligned} & v_k(G_1 + g_{kl}; H_0 = h_{kl}) - v_k(G_1; H_0) > v_i(G_1 + g_{ij}; H_0 = h_{ij}) - v_i(G_1; H_0) \\ & v_l(G_1 + g_{kl}; H_0 = h_{kl}) - v_l(G_1; H_0) > v_j(G_1 + g_{ij}; H_0 = h_{ij}) - v_j(G_1; H_0) \end{aligned}$$

Let  $v^-$  and  $v^{-0}$  denote the respective degrees in  $G_1$  of the maximally connected agents belonging to  $C = C(H_0)$  and  $C_0 = C_0(H_0)$ , and  $s^-$  and  $s^{-0}$  denote their respective friendship levels. Further, let:

$$v^-(G_1; H_0) = v^- + c_0^- s^- \tag{11}$$

$$v^{-0}(G_1; H_0) = v^{-0} + c_0^{-0} s^{-0} \tag{12}$$

Then the maximally connected agent in  $C = C(H_0)$  will propose to transform a neutral or hostile relationship with the maximally connected agent in  $C_0 = C_0(H_0)$  if:

$$v^-(G_1; H_0) + s^- > v^{-0}(G_1; H_0) + s^{-0}$$

Now, recalling that  $w(\tau + 1) > 0$  for all  $\tau \geq 0$  by assumption A.2, let us define:

$$\theta(G_1; H_0) = \min \left\{ \frac{\delta(G_1; H_0) + c^{-\tau} + 1}{w^{-\tau} + 1}, \frac{\theta(G_1; H_0) + c^{-\tau} + 1}{w(\tau + 1)} \right\}; \quad (15)$$

$\theta(G_1; H_0)$  is the threshold value of discord at which at least one agent, given their current degree in  $G_1$  and friendship in  $H_0$ , is indifferent towards transforming a link in  $H_0$ . An examination of (15) shows that, ceteris paribus, two agents have an incentive to transform their relationship if its cost  $c$  is low, the respective hostility levels they face is low (i.e., their  $\tau$ -values are high), and link formation in  $G_1$  has conferred a high enough degree on each to make it attractive to overcome any hurdle posed by their mutual discord. If  $\theta_{ij} = \theta_{ji} < \theta(G_1; H_0)$ , then at least one agent will get a negative payoff from revising their affinity relationship and will either not make such an overture (if it is the active agent) or will reject the overture (if it is the passive agent). From Lemma 1, this is also true for all pairs of agents drawn from the two clusters. Thus the existing affinity relationships in  $H_0$  between the two clusters will continue to remain hostile. Recalling (1), we have the following result:

**Proposition 4** Consider the affinity network  $H_0$  and suppose that for each pair of clusters  $C_i(H_0)$  and  $C_j(H_0)$  the divergence in their core norms satisfies:

$$\theta_j(C_i(H_0)) - \theta_i(C_j(H_0)) > \frac{1}{\theta(G_1; H_0)} - 1$$

Then  $(G = G_1; H = H_0)$  is the pairwise-stable multilayer network.

If the dissonance in core beliefs is sufficiently large between each pair of clusters, then agents within a cluster have no incentive to change their cross-cluster affinity relationships in  $H_0$ . Thus the architecture of  $H_0$  remains unchanged. Consequently, the friendship levels of agents continue to remain the same as in  $H_0$ . When link formation returns to the alliance layer, then the strategic incentives to form alliances remain the same as when link formation first started in  $G$ . Since all profitable opportunities to form alliances had already been exhausted in  $G_1$ , the architecture of the alliance network remains unchanged from  $G_1$ . Thus, all alliances that are forged continue to be within clusters and we do not observe any alliances spanning disparate clusters. Despite the high degrees of potential partners in  $G_1$ , and the accompanying economies of scale, all clusters continue to remain hostile and isolated in both the affinity and alliance layers.

Once again consider  $i \in C_i(H_0)$  and  $j \in C_j(H_0)$  who are maximally connected in  $G_1$  within their respective clusters. Now suppose  $\theta_{ij} = \theta_{ji} > \theta(G_1; H_0)$ , i.e.,

$$\theta_j(C_i(H_0)) - \theta_i(C_j(H_0)) < \frac{1}{\theta(G_1; H_0)} - 1$$

Then, because the difference in their norms is relatively small, the two agents  $i$  and  $j$  have an incentive to transform their affinity relationship. Therefore, there exists at least one agent pair in the two clusters who

will expect a change in their affinity relationship. Recall that *at most* one link between two clusters will be transformed into a positive one, since other agents in the two clusters can free ride on this “friendly” link to connect to others in the opposite cluster. Therefore, the pair transforming their link generate positive externalities for all other agents in the two clusters. Note that  $H_1 \notin H_0$  because at least one neutral or hostile link in  $H_0$  has been transformed to a friendly one. Since this transformation is predicated on the mutual profitability of an alliance in  $G$ , it follows that  $G_2 \notin G_1$ . Note an important difference now from link formation in the very first iteration on  $G_0$ . In  $G_0$ , links could only be proposed to an agent *within* the cluster; in  $G_1$ , an active agent can now propose links to agent *outside* the cluster as well who are distant friends thus leading to cross-cluster alliances in  $G_2$ . Since degrees and friendships have now changed, accounting for these changes in the expressions for (11), (12) and (15) generates a new threshold value of



Figure 5: Inter-Cluster Clique in the Alliance Network

Part (b) of Proposition 5 states that if two agents from different clusters have formed an alliance, then all agents in these two clusters with greater degree and greater friendship measure will end up interconnecting with each other. The incentives can best be explained with reference to Figure 5. Agent  $i_1$  in cluster 1 and agent  $j_1$  in cluster 2 have transformed their neutral relationship to friendly in the affinity network which is indicated by the double line connecting the two agents. Note that  $i_1 = j_1 = 1$  and  $i_1 = j_1 = 1$  prior to transforming their relationship. Suppose  $i_1$  was the active agent and  $j_1$  was the passive agent when this relationship was transformed. Therefore,  $i_1$ 's incremental payoff is:

$$[v(1) + c_0(1)] + [w_{12}(2) - c(1)] \quad (16)$$

An identical expression holds for agent  $j_1$ . Now consider agent  $i_2$  in cluster 1 who belongs to a higher friendship class than  $i_1$  and also has greater degree. Then  $i_2$  will also have a mutually profitable link with  $j_1$  in the alliance network. The incremental payoff to  $i_2$  from forming an alliance with  $j_1$  is:

$$v(3) + [w_{12}(3) - c(2)] \quad (17)$$

Agent  $i_2$  does not incur the cost since it can free ride on the friendship link between  $i_1$  and  $j_1$ .<sup>16</sup> We can now compare term-wise the incremental payoffs of  $i_1$  and  $i_2$ . Note that  $v(3) > v(1)$  from A.2(b), and  $w_{12}(3) - c(2) > w_{12}(2) - c(1)$  by virtue of A.2(a) and A.3(a). Further, since  $c_0(1) <$  from (A.4), it follows that (17) strictly exceeds (16). An identical argument establishes that  $j_1$  will reciprocate the link with  $i_2$ , and that mutually profitable links will also form between

the transformed link  $h_{i_1 j_1} = +1$  in the affinity network spurs the creation of an inter-cluster clique in the alliance network composed of  $\{i_1; i_2; j_1; j_2\}$ . More generally, suppose  $\alpha$  (respectively,  $\alpha_0$ ) be the agent with the lowest hostility level in  $C_1(H_0)$  (respectively,  $C_0(H_0)$ ) who are willing to form an alliance with each other. Then all agents in  $C_1(H_0)$  with  $\alpha$ -value exceeding  $\alpha$ , and all agents in  $C_0(H_0)$  with  $\alpha$ -value exceeding  $\alpha_0$ , will also have a mutually profitable alliance. Thus, we have an  $(\alpha; \alpha_0)$ -clique forming across two distinct clusters.

We now turn to the characterization of  $H_1$ . For each pair of clusters  $C_1(H_0)$  and  $C_0(H_0)$ , define  $\alpha_0(G; H)$  as in (15) but with respect to  $(G; H)$ . We will let  $H_1 \setminus H_0$  denote the new friendly relations that have been created and which did not exist in  $H_0$ .

**Proposition 6** Consider the limit network  $H_1$  in the pairwise-stable multilayer network  $(G; H)$ .

(a) Two clusters  $C_1(H_0)$  and  $C_0(H_0)$  are connected via a transformed friendly link in  $H_1 \setminus H_0$  if:

$$\alpha_0(C_1(H_0); C_0(H_0)) < \frac{1}{\alpha_0(G; H)} - 1 \tag{18}$$

(b) Suppose two clusters  $C_1(H_0)$  and  $C_0(H_0)$  have a mutually friendly link in  $H_1 \setminus H_0$  and that (without loss of generality)  $\alpha_0(C_1(H_0)) < \alpha_0(C_0(H_0))$ . If there exists an "intermediate" cluster  $C_{\infty}(H_0)$  such that:

$$\alpha_0(C_1(H_0)) < \alpha_0(C_{\infty}(H_0)) < \alpha_0(C_0(H_0)) \tag{19}$$

then both  $C_1(H_0)$  and  $C_0(H_0)$  have a friendly link with  $C_{\infty}(H_0)$  in  $H_1 \setminus H_0$ .

(c) Two clusters  $C_1(H_0)$  and  $C_0(H_0)$

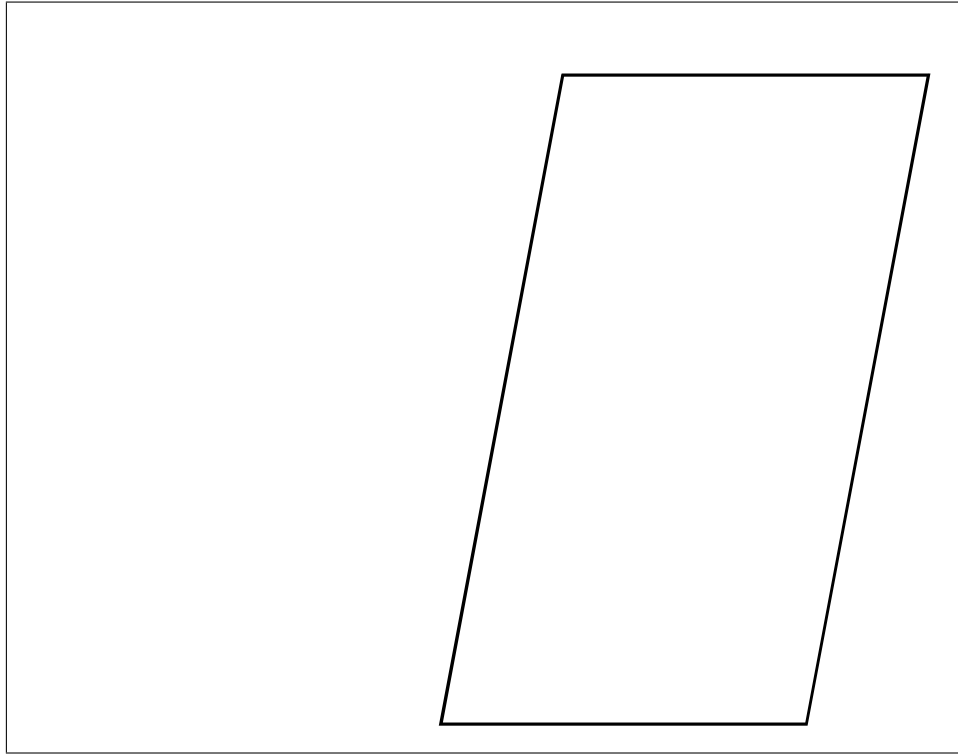


Figure 6: Overlapping Cliques in the Alliance Network

We now draw out the role of *bridge agents* who facilitate alliances across clusters. Consider Figure 6 which assumes that the relationship given by (18) holds between clusters 1 and 2, and between clusters 2 and 3. Further, cluster 2 is intermediate between the other two in the sense of (19). Finally, the relationship between clusters 1 and 3 is characterized by (20). The transformed relationship  $h_{i_2 j_2} = +1$  (shown by the double line) in the affinity network connects clusters 1 and 2 and prompts the creation of the inter-cluster clique  $f i_1; i_2; j_1; j_2 g$  with agents  $i_1$  and  $j_1$  free riding on the friendly path created between the two clusters by  $h_{i_2 j_2} = +1$ . Likewise, the transformed relationship  $h_{j_1 k} = +1$  precipitates the creation of the inter-cluster clique  $f k; j_1 g$ ; in this particular example, we are assuming that the degree and friendship measure of  $j_2$  is not sufficient for an alliance with  $k$  and the inclusion of  $j_2$  in the clique. Also, despite a friendly path now existing between clusters 1 and 3, the divergence in their core beliefs dissuades agents  $k$  and

are also in the intermediate cluster but are not members of the inter-cluster cliques due to high hostility in the affinity network and low degree in the alliance network.

The more interesting case is when, despite the relationship between clusters 1 and 3 characterized by (20), the agents in these two clusters end up forming an alliance through the aegis of agents in cluster 2 who serve as bridge agents. Due to condition (20), agents such as  $k$  and  $i_1$  are sufficiently divergent in terms of their norms such that their incremental payoff does not cover the cost of transforming their relationship. However, with the friendly path that is now created through cluster 2, these agents can eschew the transformation cost of  $\gamma$  and free ride to a mutually profitable alliance in the affinity network. Therefore, the links  $h_{i_2 j_2} = +1$  and  $h_{j_1 k} = +1$  confer positive externalities within clusters as well as across clusters permitting the formation of alliances between disparate agents who otherwise would not have an incentive to ally with each other. Therefore, through the bridge provided by agents  $j_1$  and  $j_2$ , we have an inter-cluster clique  $\{k; i_1; i_2; j_1; j_2\}$  that spans three clusters.

## 5 Motivating Examples

We now provide a set of real world examples to substantiate our main results. To illustrate the *emergence of an NSG alliance structure among agents belonging to the same cluster in the affinity network with high centrality (degree) corresponding to greater friendship measure*, we look to the Western Pacific and its overlapping relationships in  $G$  in Figure 7. These relationships run largely through the United States in an NSG type configuration. At the north end of the figure, the NATO security alliance forms a connected alliance that runs dominantly through the United States to form paths to other states and alliances. Prior to 2023, the graph can be partitioned into a set of cliques (NATO, AUKUS, FIVE EYES) and an independent set (Japan, Republic of Korea, India, Taiwan, and Micronesia) forming a partial star architecture with the United States at its center. The United States' position as an economic partner, its cultural ties, and its role as a democratic security guarantor, imparts to it the highest aggregate ranking of friendship (lowest hostility). Thus, in accordance with our result, the United States has an exceptionally high centrality in the affinity network and is by far the highest degree node in the alliance network.

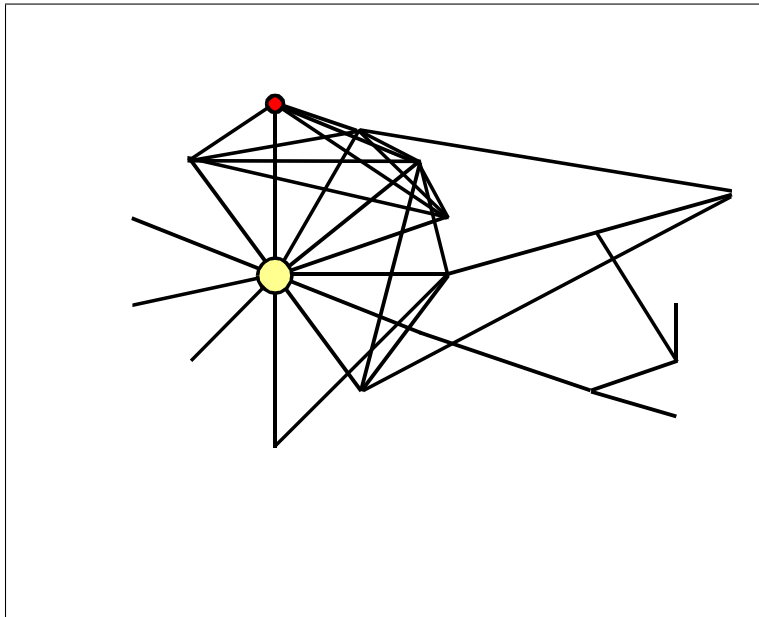


Figure 7: Alliance Network in the Pacific

Our model indicates that *when two agents belong to separate clusters in the ac nity network, then as a precursor to forming a link in the alliance network this pair of agents have to establish a bridge in the ac nity network*. In the post-pandemic years (2022-23), Japan made moves in both the H and G layers. *First*, Japan and Australia signed a security cooperation agreement (a link in G) that built on years of increasing economic ties (establishing a bridge in H). *Second*, Japan and South Korea are making significant diplomatic and economic investments at the encouragement of the United States in flipping their negative relationship into a positive one (establishing a bridge in H) as a precursor to security agreements (connecting in G). These moves are changing the existing star network in the Pacific where the United States underwrote security for all states (the graph depicted in Figure 3(a)) into a more interconnected web of alliances (Figure 3(d)).

A real world application of our result that there is *alliance formation among enemies when their disagreement over a norm is relatively small* is provided by Balkan conflict from 1992-95 (Becker et al 2023). In this case, there were three (singleton) clusters, with each cluster corresponding to an ethnic group – Bosniak, Croat and Serb. Figure 5 depicts the multilayered links between the three ethnic groups indexed by layer. The period of the conflict (1992-95) can be divided into three periods, indexed by  $t \in \{1, 2, 3\}$ , and the alliance structure that prevailed during these three periods is also indicated in Figure 5. At the multilayer network's most base *ac nity* layer, all relationships were negative since each group held deep and lasting animosities towards the others which predated the war by centuries. The other layers indicate







## 6.1 Conflict Models

We will show how the reduced form payoff function given by (4) can be deduced from an explicit conflict model. The conflict game is drawn from Baliga and Sjöström (2012). Consider agent  $i$  with neighborhood  $N_i(G)$ . Then agent  $i$  can find itself paired in pairwise conflict games – one with any of its own enemies, and one with any of the enemies of its allies. Therefore, agent  $i$  can be involved pairwise in  $|N_i(G)|$  conflict games. In each conflict game, the paired agents move simultaneously to choose either an aggressive action of *hawk* or a peaceful action of *dove*.

Recalling (3), maximizing payoffs in either the hawk-hawk or mixed strategy Nash equilibrium is equivalent to maximizing  $\pi_i(G; H_0)$  given by (4).<sup>18</sup> Therefore, our payoff specification is grounded in a proper conflict game.

Remark: (*Other parametric restrictions*) Baliga and Sjöström (2012) classify agents as *opportunistic* if hawk (dove) is a best response to dove (hawk), and *hawk-dominant* (*dove-dominant*) if hawk (dove) is a dominant strategy. The dove-dominant case once again yields an empty affinity network. The hawk-dominant case, and the pure strategy Nash equilibria of the opportunistic case yield total payoff given by (21). Thus, once again our reduced form payoff applies. The mixed strategy Nash equilibrium of the opportunistic case yields an anomalous result where agents would seek more hostilities and less allies.

## 6.2 Static Link Formation Game

We will explore a static alliance formation game adapted from Dutta et al. (1998) and show that the set of strongly stable equilibrium networks display the NSG architecture within each cluster. Thus, we provide an alternative approach to link formation that is distinct from the dynamic game as well as demonstrate the robustness of the NSG architecture. Each agent makes an announcement of intended alliances. An *announcement* by agent  $i$  is of the form  $\mathbf{s}_i = (a_{ij})_{j \in i}$ . The intended alliance  $a_{ij} \in \{0, 1\}$ , where  $a_{ij} = 1$  means that  $i$  intends to form an alliance with  $j$ , while  $a_{ij} = 0$  means that  $i$  intends no such alliance. Let  $\mathcal{S}_i$  denote the set of announcements, or strategies, of agent  $i$ . An alliance between agents  $i$  and  $j$  is formed if and only if  $a_{ij} = a_{ji} = 1$ . We denote the formed link by  $g_{ij} = 1$  and the absence of a link by  $g_{ij} = 0$ . A strategy profile  $\mathbf{s} = \{s_1, s_2, \dots, s_n\}$ , consisting of a strategy for each agent, induces a network  $G(\mathbf{s})$ . To simplify the notation we shall often omit the dependence of the network on the underlying strategy profile.

A strategy profile  $\mathbf{s} = \{s_1, s_2, \dots, s_n\}$  is *Nash* if and only if  $\pi_i(G(\mathbf{s}_i; \mathbf{s}_{-i}); H_0) \geq \pi_i(G(\mathbf{s}_i'; \mathbf{s}_{-i}); H_0)$ , for all  $\mathbf{s}_i' \in \mathcal{S}_i$  and for all  $i \in N$ , where  $\mathbf{s}_{-i}$  is the strategy profile of all agents other than  $i$ .

**Definition (Strong Stability):** A network  $G$  is said to be *strongly stable* if for any coalition  $S$  and any  $G^0$  that can be obtained from  $G$  through deviations by  $S$ ,  $u_i(G^0; H_0) > u_i(G; H_0)$  for some  $i \in S$  implies that  $u_j(G^0; H_0) < u_j(G; H_0)$  for some  $j \in S$ .

The definition of strong stability that we employ is due to Dutta and Mutuswami (1997). According to their definition, if a network  $G$  is *not* strongly stable, then there exists a coalition  $S$  that can deviate to some network  $G^0$  in which *all* members of  $S$  are strictly better off.

**Definition (Equilibrium Network):** A network  $G$  is an *equilibrium network* if there is a Nash strategy profile supporting  $G$ , and the network  $G$  is strongly stable.

In our network setting, the only *unilateral* decision that an agent has is to sever alliances. The first property of an equilibrium network is, therefore, that no agent should have an incentive to delete any subset of its alliances. Note that forming an alliance is a *bilateral* decision requiring agreement by both agents. The second property of an equilibrium network states that, for any coalition, the member agents have no incentive to bilaterally form alliances that did not exist in the equilibrium network. The second property permits a refinement of the set of Nash networks that satisfy the first property. The next result shows that all equilibrium networks display an intra-cluster NSG structure in which the neighborhood of an agent with a lower friendship measure is nested within the neighborhood of an agent with a higher friendship measure.

**Proposition 7** *An equilibrium network exists. In an equilibrium network  $G$ , all agents belonging to the same cluster form an alliance with an NSG architecture such that if  $u_i > u_j$ , then  $N_i(G) \supseteq N_j(G) \cap \{j\}$ .*

### 6.3 Non-Separable Benefits and Costs

We had assumed additively separable benefit and cost functions. This permitted us to avoid interaction between degrees of agents, or between degree and hostility. However, our results would continue to hold under a more general non-separable specification with suitable restrictions on the interaction terms. We now spell out the precise set of restrictions that are needed. Suppose the gross benefit to agent  $i$  from a link with agent  $j$  is more specified as  $b_{ij}$ . The function  $b_{ij} : Z_+^2 \rightarrow R_+$  is assumed to satisfy the following conditions:

**Assumption (A.2)\*:** For all  $i, j \in N$ :

- (a)  $b_{i+1, j} \geq b_{i, j}$



## 6.4 Endogenous Affinity Network and Norms

We have assumed that in the initial position the affinity network is given. As a first step towards a microfounded affinity network, we can assume that the affinity network is initially empty, and draw upon the definitive analysis of Hiller (2017) to augment our link formation game with the prior formation of an affinity network. Adapting Hiller, we can craft an *affinity formation game* as follows. Let  $f_1(\cdot); f_2(\cdot); \dots; f_K(\cdot)$  denote the distribution of norms that separates agents. The parameter  $\kappa$  captures the dimension on which the norm is based, i.e.,  $\kappa \in \{2, \dots, K\}$  (culture, ideology, politics, security). Assume that  $n_k$  denotes the number of agents who subscribe to the norm  $f_k(\cdot)$  such that  $n_k \geq 1$ ,  $n_k \leq n_{k+1}$ , and  $\sum_{k=1}^K n_k = N$ . Each agent is endowed with a given intrinsic level of strength that is normalized to unity. An agent can augment this strength through positive connections in the affinity network with agents *who share the same norm*. Formally, the strength gained by an agent  $i$  from establishing a positive connection with agent  $j$  in the affinity network is equal to 1 if  $f_i(\cdot) = f_j(\cdot)$  and 0 otherwise.

Each agent simultaneously proposes positive (friendship) or negative (enemy) links to other agents in the





## 8 Appendix

Proof of Theorem 1: We will show that cycles cannot emerge along an improving path in  $G$ . Let  $g_{ij} = 1$

while for agent  $j$  it follows from A.2(b), A.2(a) and A.3(a) that:

$$v_j(G^0) + g_{kj} - v_j(G^0) + w_k(G^0) + g_{kj} - c(k) \\ v_j(G + g_{ij}) - v_j(G) + [w_i(G + g_{ij}) - c(i)] = 0$$

and at least one LHS is strictly positive. Therefore, agents  $k$  and  $j$  have a mutually profitable link in  $G^0$ .

■

Proof of Proposition 2: We will first prove that  $\pi_1(H_0) \subseteq D_1(G_1)$ . Suppose to the contrary that  $i \in \pi_1(H_0) \setminus D_1(G_1)$  for  $i \geq 2$ . Thus,  $d_i(G_1) > d_j(G_1)$  for  $j \in D_1(G_1)$ . Let  $k \in N_i(G_1) \cap N_j(G_1)$  denote the agent with whom  $i$  formed a link when it had  $d_j(G_1)$  number of links, i.e, the same number of links as  $j$ . Let  $G_0(\cdot)$  denote the stage along the improving path when this link was formed, and so  $d_j(G_1) = d_i(G_0(\cdot))$ . Therefore:

$$v_i(G_0(\cdot)) + [w_k(G_0(\cdot)) + 1 - c(k)] = 0 \quad (28)$$

$$v_k(G_0(\cdot)) + [w_i(G_0(\cdot)) + 1 - c(i)] = 0 \quad (29)$$

and at least one inequality is strict. Since  $k \notin N_j(G_1)$  in the limit network  $G_1$ , it must be true that agents  $j$  and  $k$  do not have a mutually profitable link in  $G_1$ :

$$\min \{ v_k(G_1) + w_j(G_1) + 1 - c(j) ; v_j(G_1) + [w_k(G_1) + 1 - c(k)] \} < 0 \quad (30)$$

However, since  $d_j = d_i$  (given that  $i \in \pi_1$ ) and  $G_0(\cdot) \subseteq G_1$ , it follows from  $d_j(G_1) = d_i(G_0(\cdot))$  and A.3(a) that:

$$v_j(G_1) + [w_k(G_1) + 1 - c(k)] = v_i(G_0(\cdot)) + [w_k(G_0(\cdot)) + 1 - c(k)] = 0 \\ v_k(G_1) + w_j(G_1) + 1 - c(j) = v_k(G_0(\cdot)) + [w_i(G_0(\cdot)) + 1 - c(i)] = 0$$

which contradicts (30). Thus,  $d_i(G_1) = d_j(G_1)$  for all  $j \in D_1(G_1)$  and  $i \in \pi_1(H_0)$ , and hence  $\pi_1(H_0) \subseteq D_1(G_1)$ . We now prove that  $\pi_s(H_0) \subseteq D_m(G_1)$ . We have already shown in the main text that  $i_n \in \pi_s(H_0)$  and  $i_n \in D_m(G_1)$ . The same argument can be repeated for each member of

Thus  $i^{(1)} \in N_j(G)$ .

Now suppose this property is true for agents  $i^{(1)}; i^{(2)}; \dots; i^{(r)} \in C$ , i.e., these agents are the first  $r$  partners of  $i_k$  and belong to  $N_j(G) \setminus N_j(G)$ . Consider the next partner  $i^{(r+1)}$  of agent  $i_k$  and suppose this link was formed in stage  $G_0^{(r)}$  along the improving path. Suppose  $i_k$  was the active player when this link was formed. Since  $|N_j(G_0^{(r)})| = r$ , it follows from Proposition 1 that  $i^{(r+1)}$  and  $i_j$  also have a mutually profitable link when  $i_j$  is the active agent. Now suppose  $i^{(r+1)}$  was the active agent when the link with  $i_k$  was formed. Then, similar to the reasoning with  $i^{(1)}$ , agent  $i^{(r+1)}$  would have first formed this link with  $i_j$ . Therefore,  $i^{(r+1)} \in N_j(G)$ . This completes the induction step and proves the nestedness property. ■

**Proof of Theorem 2:** In the augmented link formation game, as each iteration of link formation occurs in the alliance network  $G_r, r \geq 1$ , potentially new alliances are added but none of the existing links are removed. Let  $d(G) = (d_1(G), d_2(G), \dots, d_N(G))$  denote the *degree distribution* of  $G$ .

The verification for agents  $j$  and  $l$  is identical. ■

Proof of Proposition 4: The proof is provided in the main text. ■

Proof of Proposition 5:

a. We have proved in Proposition 3 that  $G_1$  has an NSG architecture in each cluster. Now suppose this is true for  $G_r, r \geq 2$ . We will prove it for  $G_{r+1}$  by contradiction. Suppose there exists a cluster  $C \subseteq (H_0)$  with agents  $i$  and  $j$  such that  $d_i(G_{r+1}) < d_j(G_{r+1})$  but  $N_i(G_{r+1}) \neq N_j(G_{r+1})$ . In particular, there exists an agent  $k \in C \subseteq (H_0)$  such that  $k \in N_i(G_{r+1}) \setminus N_j(G_{r+1})$ . Since  $C \subseteq (H_0)$  has an NSG structure in  $G_r$ , and  $G_r \subseteq G_{r+1}$ , the link  $g_{ik} = 1$  must have been added when link formation was occurring in  $G_{r+1}$ . Thus,  $d_i(G_r) < d_j(G_r)$ . Recalling Proposition 3 which demonstrated that degree is positively correlated with friendship, it follows that  $d_i < d_j$ . Now suppose the network is  $G_{r+1}(\cdot)$  when the link  $g_{ik} = 1$  is formed in  $G_{r+1}$ . There are two possible cases.

Case I: Suppose  $i$  was the active agent and  $k$  acquiesced as the passive agent. Then, in some subsequent state  $(G_{r+1}(\cdot^0); H_r; k)$ , i.e., when  $k$  is the active agent, then  $k$  will have a mutually profitable link with  $j$ .

$$v_k(G_{r+1}(\cdot^0)) + g_{kj}; H_r \quad v_k(G_{r+1}(\cdot^0); H_r) = v_k(G_{r+1}(\cdot^0)) + w_j(G_{r+1}(\cdot^0)) + 1 - c(j)$$

From A.2(b):

$$v_k(G_{r+1}(\cdot^0)) = v(k(G_{r+1}(\cdot^0)))$$

and, since  $d_i(G_r) < d_j(G_{r+1}(\cdot^0))$ , from A.2(a) and A.3(a):

$$w_j(G_{r+1}(\cdot^0)) + 1 - c(j) > w(d_i(G_r) + 1) - c(i)$$

Therefore:

$$v_k(G_{r+1}(\cdot^0)) + g_{kj}; H_r \quad v_k(G_{r+1}(\cdot^0); H_r) > v(k(G_{r+1}(\cdot^0))) + w(d_i(G_r) + 1) - c(i) > 0$$

where the second strict inequality follows from the fact that agent  $k$  had acquiesced to a link with  $i$  when the network was  $G_{r+1}(\cdot)$ . Agent  $j$  will reciprocate because:

$$\begin{aligned} v_j(G_{r+1}(\cdot^0)) + g_{kj}; H_r \quad v_j(G_{r+1}(\cdot^0); H_r) &= v_j(G_{r+1}(\cdot^0)) + w_k(G_{r+1}(\cdot^0)) + 1 - c(k) \\ &> v(d_i(G_r)) + [w(k(G_{r+1}(\cdot^0))) + 1 - c(k)] > 0 \end{aligned}$$

where the last strict inequality follows since  $i$  had proposed a link to  $k$  in  $G_{r+1}(\cdot)$ . Therefore, it cannot be the case that when all profitable opportunities have been exhausted in  $G_{r+1}$  then agents  $k$  and  $j$  will remain unlinked.

Case II: Suppose  $k$  was the active agent when the network was  $G_{r+1}(\cdot)$ . Then, according to the link

formation protocol,  $k$  would have proposed a link with agent  $j$  rather than  $i$  because:

$$\begin{aligned}
 v_k(G_{r+1}(\cdot) + g_{kj}; H_r) - v_k(G_{r+1}(\cdot); H_r) &= v_k(G_{r+1}(\cdot)) + w_j(G_{r+1}(\cdot)) + 1 - c(j) \\
 &> v_k(G_{r+1}(\cdot)) + [w_i(G_r) + 1 - c(i)]
 \end{aligned}$$

and once again agent  $j$  will accept the proposal. Therefore, once again we have a contradiction.

It follows that each intra-cluster architecture in  $G_{r+1}$  will have an NSG architecture. Since  $G$  is reached in a finite number of steps, it follows that the result also holds for  $G$ .

b. Let  $i \in C_1(H_0) \setminus C_1(H)$  and  $j \in C_0(H_0) \setminus C_0(H)$ . Let  $H_{i,j}$  denote the affinity network  $H$  in which there is no friendly link between clusters  $C_1(H_0)$  and  $C_0(H_0)$ . There are two possible cases:

Case I: Suppose  $i$  and  $j$  incurred the cost of transforming their affinity relationship allowing all other agents in the two clusters to free ride on the friendly path they have created. Following the same argument as Lemma 1, for any  $k_1 \in C_1(H_0) \cap i^c$  and  $k_2 \in C_0(H_0) \cap j^c$ :

$$\begin{aligned}
 v_{k_1}(G; H_{i,j}) - v_{k_1}(G; H) &> 0 & v_{k_1}(G; H_{i,j}) - v_{k_1}(G; H) &> 0 \\
 v_{k_2}(G; H_{i,j}) - v_{k_2}(G; H) &> 0 & v_{k_2}(G; H_{i,j}) - v_{k_2}(G; H) &> 0
 \end{aligned}$$

where at least one of the last inequality in each case is strictly positive. Since  $k_1$  and  $k_2$  free ride, it follows that:

$$\begin{aligned}
 v_{k_1}(G; H) - v_{k_1}(G; H_{i,j}) &> 0 & v_{k_1}(G; H) - v_{k_1}(G; H_{i,j}) &> 0 \\
 v_{k_2}(G; H) - v_{k_2}(G; H_{i,j}) &> 0 & v_{k_2}(G; H) - v_{k_2}(G; H_{i,j}) &> 0
 \end{aligned}$$

and the result follows.

Case II: Suppose a pair of agents, where at least one agent differs from  $i$  or  $j$ , were the ones transforming their affinity relationship. Call this pair of agents transforming their affinity relationship as  $\varphi \in C_1(H_0) \setminus C_1(H)$  and  $\psi \in C_0(H_0) \setminus C_0(H)$ , where  $s \neq i$  and  $s \neq j$ . Following the same argument as that in Case I, all agents with friendship measures greater than or equal to those of  $\varphi$  and  $\psi$  will also have an incentive to form an alliance. Now consider agents  $i$  and  $j$  from the statement of the proposition. These two agents will free ride on the friendly link created by  $\varphi$  and  $\psi$  and have a profitable alliance by hypothesis. Thus, for any two agents  $k_1$  and  $k_2$  whose friendship measures are greater than or equal to those of  $i$  and  $j$  and who also free ride, we have:

$$\begin{aligned}
 v_{k_1}(G; H) - v_{k_1}(G; H_{i,j}) &> 0 & v_i(G; H) - v_i(G; H_{i,j}) &> 0 \\
 v_{k_2}(G; H) - v_{k_2}(G; H_{i,j}) &> 0 & v_j(G; H) - v_j(G; H_{i,j}) &> 0
 \end{aligned}$$

Proof of Proposition 6: The proof follows from the definition of the threshold value,  $\theta_i(G; H)$ . ■

Proof of Proposition 7: To save space, we will suppress reference to  $H_0$ .

(Existence): We first establish existence. Recall that all alliances are formed within clusters. Consider the network in which each cluster  $C$  is complete, i.e., all agents in each cluster are mutually interconnected. Denote this network as  $G^c$ . If it is an equilibrium, then we are done. Otherwise, there exists a coalition  $S^0$  and a network  $G^0$  that can be obtained from  $G^c$  by  $S^0$  such that  $v_i(G^0) > v_i(G^c)$  for all  $i \in S^0$ . Since all alliances are intra-cluster, it implies that  $S^0 \subset C$  for some cluster  $C$ . Specifically:

$$v_i(G^0) = v_i(G^c) - c_0(i) + \sum_{j \in N_i(G^0)} w_j(G^0) - c(j) > v_i(G^c); \quad i \in S^0$$

Since no new links could be added in  $G^c$ , the deviation must involve members in  $S^0$  deleting their links. This implies in particular that in the cluster  $C$ :

$$v_j(C) - 1 + [w_j(C) - c(j)] < 0; \quad i \in S^0; \quad j \in N_i(G^c) \cap N_i(G^0) \quad (31)$$

If  $G^0$  is an equilibrium, then we are done. Otherwise, there exists a coalition  $S^{00}$  that can obtain a network  $G^{00}$  in which each member is strictly better off. We claim that this movement from  $G^0$  to  $G^{00}$  can only involve a deletion of links. Suppose to the contrary that the movement from  $G^0$  to  $G^{00}$  involves addition of links and let  $S^0 \setminus S^{00}$  denote the non-empty subset of agents who are involved in forming alliances, either among themselves or with others in  $S^{00} \cap S^0$  in the move from  $G^0$  to  $G^{00}$ . Note that this intersection cannot be empty because firms in  $C \cap S^0$  are completely connected among themselves; thus a member of  $S^0$  has to be involved if new links are created starting from  $G^0$ . Consider any  $i \in S^0 \setminus S^{00}$ . Since  $i$  was completely connected in  $G^c$ , and deleted links in the move to  $G^0$ , any new alliance that it forms in the move to  $G^{00}$  must be with some agent  $j \in N_i(G^c) \cap N_i(G^0)$  with whom it earlier dissolved an alliance. Since the deviation to  $G^{00}$  is strictly profitable:

$$v_i(G^{00}) - 1 + \sum_{j \in N_i(G^{00})} w_j(G^{00}) - c(j) > 0 \quad (32)$$

However,  $v_j(G^{00}) - 1 - \sum_{i \in N_j(G^{00})} w_i(G^{00}) + c(j) < 0$  and  $j \in N_i(G^c) \cap N_i(G^0) \cap N_j(G^{00})$

in the move to  $G^{00}$







and  $j$  respectively from forming a link in  $G^0$  is equal to:

$$\sum_{i \in N_i(G^0)} \binom{0}{k; i} + \binom{0}{k+1; j+1} c_{j; j+1} \quad (35)$$

$$\sum_{j \in N_j(G^0)} \binom{0}{j; i} + \binom{0}{j+1; k+1} c_{k; k+1} \quad (36)$$

Note that  $N_i(G) = N_k(G^0)$  and  $N_j(G) = N_j(G^0)$ . For all  $i \in N_i(G)$ , since  $\binom{0}{k; i} = i$ , it follows from respectively parts (b) and (d) of A.2 that  $\binom{0}{k; i} = \binom{0}{i; j}$  and  $\binom{0}{i; j} = \binom{0}{i; j}$ . Further, from A.2 (a),  $\binom{0}{k+1; j+1} = \binom{0}{i+1; j+1} = \binom{0}{i+1; j+1}$ . Further, since  $\binom{0}{j; j} = j$ , from A.3 (a),  $c_{j; j+1} = c_{j; j+1}$ . Therefore, each term in (35) dominates the corresponding term in (33). Likewise, noting that  $\binom{0}{k; i} = i$ , each term in (36) dominates the corresponding term in (34). This proves the result. ■

**Proof of Lemma 1 for the Non-separable Case:** Dropping reference to  $G_1$ , we will let  $i = i(G_1)$  and  $i+1 = i(G_1 + g_{ij})$ . Then,  $i(G_1 + g_{ij}; H_0) = h_{ij} = i(G_1; H_0)$  is equal to:

$$\sum_{h \in N_i(G_1)} \binom{0}{i; h} + [c_{i; i} - c_{i+1; i+1}] + \sum_{j \in N_j(G_1)} \binom{0}{i+1; j+1} c_{j; j+1}$$

Similarly,  $k(G_1 + g_{kl}; H_0) = h_{kl} = k(G_1; H_0)$  is equal to:

$$\sum_{h \in N_k(G_1)} \binom{0}{k; h} + [c_{k; k} - c_{k+1; k+1}] + \sum_{l \in N_l(G_1)} \binom{0}{k+1; l+1} c_{l; l+1}$$

Since there is an NSG structure within each cluster with degree positively related to friendship,  $N_k(G_1) \subseteq N_i(G_1)$  and thus  $i \geq k$ . It follows from parts (b) and (c) respectively of A.2 that:

$$a) \sum_{i \in N_i(G_1)} \binom{0}{i; i} = \sum_{i \in N_i(G_1)} i = \sum_{i \in N_i(G_1)} i$$

and

terms:

$$c(k; i) - c(k+1; i+1) - c(k; k) - c(k+1; k+1) \quad (38)$$

From (37) and (38), it follows that:

$$c(i; j) - c(i+1; j+1) - c(k; k) - c(k+1; k+1)$$

Finally, note that  $kl = ij$ ,  $j < i$  and  $j < k$ . Therefore, from A.2 (a):

$$c(i+1; j+1) - c(k+1; j+1) - (c(k+1; i+1))$$

and from A.3 (a):

$$c(j; j+1) - c(i; j+1) - c(i; i+1)$$

It follows that:

$$i(G_1 + g_{ij}; H_0 - h_{ij}) - i(G_1; H_0) - k(G_1 + g_{kl}; H_0 - h_{kl}) - k(G_1; H_0)$$

The verification for agents  $j$  and  $l$  is identical. ■

## References

- [1] R. Aumann and R. Myerson, 1988. Endogenous Formation of Links between Players and Coalitions: An Application of the Shapley Value. In *The Shapley Value*, (A. Roth, editor), 175-191, Cambridge University Press.
- [2] S. Baliga and T. Sjöström, 2012. The Strategy of Manipulating Conflict, *American Economic Review*, 102(6), 2897-2922.
- [3] Becker C., P. Devine, H. Dogo and E. Margolin, 2022. Marking Territory: Modeling the Spread of Ethnic Conflict in Bosnia and Herzegovina, 1992-1995 SSRN: <https://ssrn.com/abstract=3160015>.
- [4] F. Bloch, 2012. Endogenous Formation of Alliances in Conflict. In *The Oxford Handbook of the Economics of Peace and Conflict* (M.R. Garaskel and S. Skaperdas, editors), 473–502, Oxford University Press.
- [5] D. Cartwright and F. Harary, 1956. Structural Balance: A Generalization of Heider's Theory, *Psychological Review*, 63(5), 277–293.
- [6] J.A. Davis, 1967. Clustering and Structural Balance in Graphs, *Human Relations*, 20(2), 181–187.
- [7] B. Dutta, A. van den Nouweland and S. Tijs, 1998. Link Formation in Cooperative Situations, *International Journal of Game Theory*, 27, 245–256.

[8] B. Dutta and S. Mutuswami, 1997. Stable networks, *Journal of Economic Theory* 76, 322–344.