

An application of formal argumentation: Fusing Bayesian networks in multi-agent systems

Søren Holbech Nielsen ^{a,*}, Simon Parsons ^b

^a *Department of Computer Science, Aalborg University, Aalborg, Denmark*

^b *Department of Computer and Information Science, Brooklyn College, City University of New York, Brooklyn, 11210 NY, USA*

Received 19 October 2006; received in revised form 28 March 2007; accepted 16 April 2007

Available online 29 April 2007

Abstract

We consider a multi-agent system where each agent is equipped with a Bayesian network, and present an open framework for the agents to agree on a possible consensus network. The framework builds on formal argumentation, and unlike previous solutions on graphical consensus belief, it is sufficiently general to allow for a wide range of possible agreements to be identified.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Argument in agent systems; Argumentation frameworks; Application; Bayesian networks

1. Introduction

Lately research in distributed systems has intensified, spurred by increased availability of sophisticated electronic devices and cheap networking equipment. Within this field, the crossover area of multi-agent systems (MAS), that incorporates parts of artificial intelligence research, has been heavily researched, with each device being modeled as an autonomous agent capable of acting on its environment, reflecting on observations of its surroundings and communicating with other agents. To implement such reflective capabilities, a model-based agent architecture is often employed, where the agent carries within it a formal model of its surroundings and acts on mathematical inferences drawn from this model. Hence, the more accurate the model is the more successful the agent will be in achieving its objectives. Therefore it is beneficial for the agent to (i) update its model when observations of the agent's surrounding indicate that the model is inaccurate, and (ii) communicate with other agents about their models and alter its own model to reflect model aspects common to the models of most other agents.

In this text we investigate how Bayesian networks (BNs) can be used as internal models in a multi-agent setting, and more specifically how (ii) can be implemented with the help of argumentation theory. Previously the two methodologies have mainly been studied together with a view to incorporating the efficiency and precision of BNs into argumentation theory (e.g. [24]), or as an exercise in converting models of one theory into models of the other (e.g. [27] and [30]). Here, we instead try to exploit strong points of both reasoning methods: BNs constitute a com-

* Corresponding author.

E-mail addresses: soeren.holbech@gmail.com (S.H. Nielsen), parsons@sci.brooklyn.cuny.edu (S. Parsons).

pact, elegant, and mathematically correct framework for drawing diagnostic or causal inferences from observations to hypothesis variables, even in face of noisy observations, and they can be constructed and altered automatically from observation data, and are thus good choices for internal models. Argumentation theory, on the other hand, provides a methodology for transparently extracting a consistent “truth” from a set of conflicting and/or overlapping views, and furthermore has strong roots in dialectics, which makes distributed agent-oriented implementations natural.

Playing along with these strengths, we envision a MAS of cooperating agents, where each agent has a BN as a model of the domain it is situated in, and aim at providing a framework built on principles of formal argumentation theory in which the agents, starting from their individual domain models, can end up agreeing on a single network representing their joint domain knowledge.

The task of fusing several BNs into one compromise BN has previously been addressed with an a priori specified view to what constitutes a compromise [4,10,14,21,23], with no apparent consensus on the goal of network fusion among the authors, and has mainly been considered a centralized operation. Matzkevich and Abramson [14] disregard the strength of independency statements in the input networks, and seek to obtain a graph that contains all arcs from the input networks or their reverses. Sagrado and Moral [4] similarly disregard the quantitative parts of the input networks, and gives rules for constructing graphs that imply either all independency statements implied by at least one of the input BNs, or only those independency statements implied by all input BNs. Richardson and Domingos [23] also disregard the quantitative parts of the input networks, and construct a prior distribution over all graph structures from the input BNs. This prior is then used as basis for a standard greedy search based on a separate database. Finally, Pennock and Wellman [21] derive a series of impossibility results for the general problem of combining probability distributions, and Il and Chajewska [10] uses some of the results from [21] to adapt standard greedy search score functions to use the input BNs, rather than a database, as their basis. The differing objectives of these papers stems from differences in interpretations of BNs: They are seen as either specifying flow of information, specifying independencies between variables, representing expert experience, or summarizing data. In this paper, we do not commit ourselves to any specific compromise objective. Rather, we establish a general framework in which any kind of compromise can be reached, with the exact nature of this specified by a compromise score function and possibly a heuristic for walking search trees. The advantages of our approach include that a general purpose argumentation engine can be implemented, and reused in contexts with different definitions of compromise; efficient distributed implementations are natural; in cases where agents almost agree a priori, little information need to be shared among the agents; and anytime compromises can be achieved.

The text is structured as follows: In Section 2 we briefly cover the methodology that we need. In particular, we need some concepts of graph theory, the main theory of BNs and their equivalence classes, and finally theory of an argumentation framework. Following this, in Section 3 we specify the problem that we wish to solve more formally, and present the strategy for doing so. The actual results then follow. First, Section 4 describes how we encode BNs in our framework. Then Section 5 presents a formal argumentation system whose preferred extensions correspond to proper BNs, and Section 6 contains the debating guide-lines that must be followed by agents. We end with a discussion of these results in Section 7.

2. Preliminaries

Before we proceed with the theory needed for the results later in the paper, we briefly clarify the notation: as a general rule, sets are printed in boldface type (\mathbf{S} , $\mathbf{pa}(V)$, ...), and individual entities are printed in plain letters (A , x_i , ...), except data structures, which are printed in calligraphic letters (\mathcal{G} , \mathcal{A} , ...), and general classes, which are printed in gothic type (\mathfrak{C}). By convention we use \equiv to mean “is defined to be” or “defined as”, and decline to use { and } when listing singleton sets. The material is necessarily heavy on definitions, and to help the reader refer back, each term is *emphasized* where it is defined.

2.1. Graphs

Here we briefly introduce the terms of graph theory and notation that are used in the remainder of the text. For more elaboration on the concepts introduced, see e.g. [13].

A *graph* is a pair $\mathcal{G} \equiv (X, E)$, where X is the finite set of *nodes* of the graph, and $E \subseteq (X \times X) \setminus \{(X, X) : X \in X\}$ is a set of ordered pairs called *edges* of the graph. For any two nodes X and Y , if both (X, Y) and (Y, X) is in E we

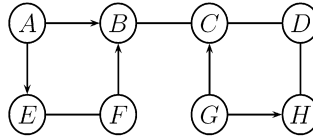


Fig. 1. A simple graph.

say that there is an *undirected link* (or just a *link*) between X and Y and that they are *neighbors*. If only (X, Y) is in E , then we say that there is an *arc from* X to Y , that X is a *parent* of Y , and that Y is a *child* of X . The set of parents of a node X is denoted $\mathbf{pa}(X)$. If X is a neighbor, parent, or child of Y , we say that X and Y are *adjacent*. A triple of nodes (X, Y, Z) is said to constitute a *v-structure*, if Y is a child of both X and Z , and X and Z are not adjacent. When depicting a graph we use circles for nodes, lines for links, and arrows for arcs, as shown in Fig. 1. From the figure it can be seen that $\mathbf{pa}(C) = G$, that C is a neighbor of both B and D , that C has no children, and that (A, B, F) is the only v-structure in the graph.

For two nodes X_1 and X_k , we say that there is a *path* from X_1 to X_k , if either (X_1, X_k) is in E , or there are distinct nodes X_2, \dots, X_{k-1} different from X_1 and X_k , such that (X_{i-1}, X_i) is in E for all $2 \leq i \leq k$. We denote the path (X_1, \dots, X_k) and say that its *length* is $k - 1$. If, for one of the edges (X, Y) in a path, (Y, X) is not in E then we say that the path is *directed*, if not it is *undirected*. A path (X, \dots, Y) , where X and Y are the same node, we call a *cycle*. Given a cycle $(X_1, \dots, X_{k-1}, X_1)$ of length $k \geq 4$, we say that the cycle has a *chord* if there is a pair of nodes X_i and X_j , where $|i - j| \geq 2$ and X_i and X_j are not the pair X_1 and X_{k-1} , such that X_i and X_j are adjacent. A cycle of length $k \geq 4$ with no chord is called *chordless*. If there is a directed path from a node X to a node Y , then we say that Y is a descendant of X . The set of descendants of X is denoted by $\mathbf{desc}(X)$. If, for any two nodes X and Y in a set Y , either X is the same node as Y , or there is an undirected path from X to Y , and this holds for no other set $Z \supset Y$, then we call Y a *chain component*. For examples, consider the graph in Fig. 1: There are no cycles in the graph, $\mathbf{desc}(A) = \{B, C, D, E, F, H\}$, and $\{B, C, D, H\}$ is a chain component.

Given a graph $\mathcal{G} \equiv (X, E)$, we introduce some auxiliary definitions. Let $Y \subseteq X$, then we say that $\mathcal{G}_Y \equiv (Y, E \cap (Y \times Y))$ is the *subgraph induced* by Y . If all the edges in \mathcal{G} are arcs, we say that \mathcal{G} is *directed*. If all the edges are links, we say that \mathcal{G} is *undirected*. A graph, which is neither directed nor undirected, is a *partially directed graph*. An undirected graph containing no chordless cycles is called *decomposable*. A graph with no directed cycles is called a *chain graph*. A directed chain graph is also called an *acyclic directed graph*, traditionally abbreviated *DAG*. The graph $\mathcal{G}^u \equiv (X, \{(X, Y): (X, Y) \in E \text{ or } (Y, X) \in E\})$, obtained by replacing each arc in \mathcal{G} with a link, is called the *skeleton* of \mathcal{G} .

2.2. Bayesian networks

A *Bayesian network* (BN), over a set of discrete¹ domain variables V , is a pair (\mathcal{G}, Φ) , where $\mathcal{G} \equiv (V, E)$ is a DAG, and Φ is a set of conditional probability distributions:

$$\Phi \equiv \{P(V | \mathbf{pa}(V)): V \in V\}.$$

In the sequel we assume V to be fixed and implicitly constitute the basis for the nodes of all graphs. For more information on BNs in general, see [11].

A BN provides a probabilistic model P of the domain V , such that for any configuration $\mathbf{v} \equiv (v_1, \dots, v_n)$ over V ,

$$P(\mathbf{v}) \equiv \prod_{V \in V} P(V_i = v_i | \mathbf{pa}(V_i) = \mathbf{v}_{\mathbf{pa}(V_i)}),$$

with $\mathbf{v}_{\mathbf{pa}(V_i)}$ denoting the entries in the configuration \mathbf{v} that corresponds to the variables in the set $\mathbf{pa}(V_i)$. For any BN, the limitations on P imposed by the structure of \mathcal{G} are summarized as a set of independence constraints:

$$\{V \perp\!\!\!\perp V \setminus (\mathbf{desc}(V) \cup \mathbf{pa}(V) \cup V) | \mathbf{pa}(V): V \in V\},$$

¹ “Discrete” meaning that a variable can be in one of only a finite number of states. Other equivalent terms include “cardinal” and “ordinal”.

where the notation $X \perp\!\!\!\perp Y \mid Z$, for three disjoint sets of variables X , Y , and Z , means that the variables in X are probabilistically conditionally independent of those in Y given those in Z , i.e.

$$P(X \cup Y \mid Z) = P(X \mid Z) \cdot P(Y \mid Z),$$

whenever $P(Z) > 0$. In other words, the set of probability models that can be associated to some DAG to constitute a BN, all need to obey the independence constraints dictated by the structure [20].

Two distinct DAGs may give rise to the same independence constraints as can be seen from the two graphs in Figs. 2(a) and 2(b): They imply the independence statements $A \perp\!\!\!\perp C \mid B$ and $C \perp\!\!\!\perp A \mid B$, respectively, which by the commutativity of standard multiplication are the same. We can thus define an equivalence relation over DAGs, such that \mathcal{G} and \mathcal{H} are equivalent iff the independence constraints implied by \mathcal{G} are exactly those implied by \mathcal{H} . This leads to an important observation: when only the structure of BNs need to be determined (e.g. in situations where only independencies among parts need consideration) several DAGs can be correct answers. Therefore, identifying a single DAG, rather than the equivalence class containing it, can be a waste of effort and maybe even impossible, if only independence information is available. Hence, it can be fruitful to consider equivalence classes rather than DAGs.

Verma and Pearl [26] showed that any two DAGs imply the same independence constraints if and only if (iff) they contain the same v-structures and have the same skeleton. Thanks to this result, we can represent the equivalence class of a DAG \mathcal{G} as a graph having the same skeleton as \mathcal{G} , the same v-structures, and all edges not participating in defining a v-structure being undirected links. Clearly this graph is uniquely determined, and following Verma and Pearl [26] we shall call it the *pattern* of the equivalence class. As an example, the pattern of the equivalence class of the DAG in Fig. 3(a) is shown in Fig. 3(b).

Once a pattern has been determined for a DAG all the members of its equivalence class can be constructed by exchanging links for arcs in all possible ways under the constraints that

- (1) the resulting graphs contains no directed cycles, and
- (2) no other v-structures than those of the pattern emerge.

These two constraints sometimes cause one or more of the arcs replacing links in the pattern to have the same orientation in the resulting DAGs, no matter in what order or direction other links are converted. For instance, the link between C and E in the pattern in Fig. 3(b) cannot be exchanged with an arc going from E to C , as that would create a new v-structure (E, C, F) . Hence, it must be exchanged for an arc from C to E in all DAGs created from that pattern. The arcs arising from replacing such links, and arcs participating in defining the v-structures in the pattern, we collectively call *compelled arcs*. The graph that results from replacing links for compelled arcs wherever possible is called the *completed pattern*. Notice that by definition this graph must necessarily be unique for each equivalence class. The completed pattern of the DAG in Fig. 3(b) is shown in Fig. 4(a).

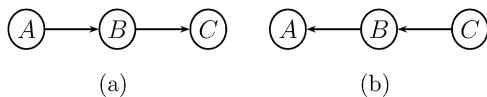


Fig. 2. Two BNs with the same independence properties.

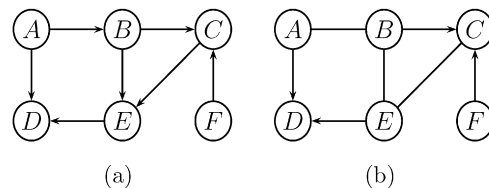


Fig. 3. (a) A graph and (b) its pattern.

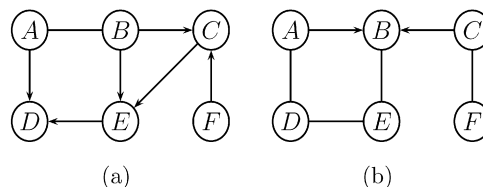


Fig. 4. The completed pattern of the graph in Fig. 3 and a graph with no consistent extension.

Now, consider the graph in Fig. 4(b). No matter how we try to exchange the links in the graph with arcs we end up either creating a directed cycle or a v-structure not present in the original graph. Thus, this graph cannot be a pattern for any equivalence class of DAGs nor be a result of extending a pattern while respecting Constraints 1 and 2, and we are therefore dealing with two classes of graphs: Those graphs that can be turned into at least one DAG, by exchanging links for arcs while respecting Constraints 1 and 2, we say *admit a consistent extension*, whereas those that cannot *do not admit a consistent extension* (terminology originally introduced by Chickering [3]). We denote the class of graphs admitting a consistent extension as \mathcal{C} . As a DAG has itself as a consistent extension, it follows that any DAG $\mathcal{G} \equiv (V, E)$ is a member of \mathcal{C} . Furthermore, the pattern $\mathcal{G}^p \equiv (V, E^p)$ of \mathcal{G} must necessarily be a member of \mathcal{C} , as must any graph $\mathcal{G}' \equiv (V, E')$, where $E \subseteq E' \subseteq E^p$, including the completed pattern of \mathcal{G} . We shall also distinguish completed patterns from other members of \mathcal{C} , i.e. $\mathcal{C}^{cp} \subset \mathcal{C}$ is the class of graphs that are completed patterns of at least one DAG.

Unfortunately, we are unaware of any simple characterization of members of \mathcal{C} . That is, potentially trying out all possible directions of undirected links in the hope of obtaining a consistent extension is not simple, but a computationally intensive procedure calling for a lot of back tracking, and it requires global investigation of the graph rather than a series of local investigations. However, Andersson et al. [1] provide the following result on members of \mathcal{C}^{cp} :

Theorem 1. A graph \mathcal{G} is a member of \mathcal{C}^{cp} iff

- (1) \mathcal{G} is a chain graph,
- (2) \mathcal{G}_C is decomposable for all chain components C of \mathcal{G} ,
- (3) if X is a parent of Y , and Z is a neighbor of Y , then X and Z are adjacent (i.e. the configuration in Fig. 5(a) does not exist as a subgraph of \mathcal{G}), and
- (4) each arc in \mathcal{G} is strongly protected (defined below).

An arc from a node X to a node Y in a graph \mathcal{G} is strongly protected if either

- there is a node Z that is a parent of X and not adjacent to Y (see Fig. 5(b)),
- there is a node Z that is a parent of Y and not adjacent to X (see Fig. 5(c)),
- there is a node Z that is a parent of Y and a child of X (see Fig. 5(d)), or
- there are non-adjacent nodes Z and W , such that both Z and W are parents of Y and neighbors to X (see Fig. 5(e)).

The bullets of the theorem all have pretty intuitive justifications. In short, items 1, 2, and 3 ensure that all compelled arcs are arcs, whereas item 4 ensures that nothing but compelled arcs are arcs. In particular, item 1 calls for the graph to be a chain graph, which is reasonable, since any graph with a fully directed cycle cannot be extended to a DAG, and any graph with a directed cycle still containing undirected links cannot be extended to a DAG without creating a new v-structure, unless the cycle is of length 3, in which case the arc(s) participating in the cycle cannot be strongly protected, and hence would be forbidden by item 4.

Item 2 is reasonable, since any graph with a chordless undirected cycle cannot be extended into a DAG without creating a new v-structure.

A graph failing to satisfy item 3 cannot be extended into a DAG without creating a new v-structure, unless the undirected link between Y and Z is directed away from Y . But then this arc is compelled, and it should have been an arc if the graph was to be a completed pattern.

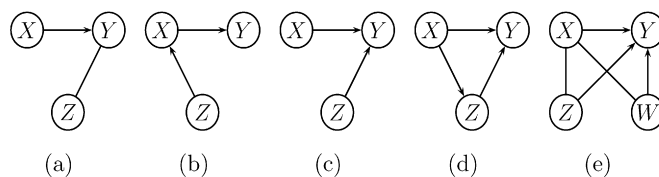


Fig. 5. Configurations described in Theorem 1.

As for item 4, consider that for any arc in a completed pattern, it holds by definition that this arc cannot be reversed without changing the set of v-structures or introducing a directed cycle. Going through each case of the definition of “strongly protected”, it is easy to see that changing the direction of these arcs does result in a cycle, a new v-structure, or the destruction of a v-structure.

Clearly, both items 3 and 4 are easily checked locally, and items 1 and 2 can be checked without any backtracking. Hence, Theorem 1 provides an efficient means for identifying members of \mathcal{C}^{cp} and thus exactly those graphs that are of interest, when only the structure of a BN needs to be identified.

2.3. Argumentation systems

Argumentation is a relatively new approach to extracting consistent knowledge from a possibly inconsistent knowledge base while possibly respecting some assumptions on the relative epistemological worth of statements in the knowledge base. The approach can be seen as an alternative to formal non-monotonic and defeasible logics. As the research area of argumentation systems is still fairly new, no single methodology has yet to stand out as the main approach, and a lot of problems are still not solved. However, the argumentation approach shows great potential as a general purpose reasoning mechanism, which is why we use it here.

Due to the lack of an established framework for analysis of argument systems, it has been necessary to pick one from a large pool of these (e.g. [8,12,22,25,29]). The framework we have picked for our purpose is the framework of [19], which generalizes that of [7], as this is an abstract framework, which leaves the underlying language and reasoning unspecified.

An *argumentation system* is a pair $\mathcal{A} \equiv (A, \triangleright)$, where A is a set of *arguments*, and $\triangleright \subseteq (2^A \setminus \{\emptyset\}) \times A$ is an *attack relation*.² The exact nature of an argument is left unspecified, but examples built on a language akin to propositional logic could be:

- “The sun is shining, so it is not raining”,
- “I saw a pigeon yesterday, so it is not raining”, and
- “It is not raining”.

For two sets of arguments $S \subseteq A$ and $S' \subseteq S$ and an argument A , if $S' \triangleright A$ then S is said to *attack* A , and A is said to *be attacked by* S . If no proper subset of S' attacks A , then S' is called a *minimal attack* on A . An example of an attack that would intuitively make sense is

“The sun is shining” \triangleright “It is raining”.

A *semantics* of an argumentation framework is a definition of which arguments in the framework that should be accepted by a rational individual. [7] and [19] work with a wide range of semantics, but we only introduce those we need here. First of all, we define a set of arguments $S \subseteq A$ as being *conflict-free*, if there is no argument A in S such that S attacks A . Intuitively, claiming a conflict-free set of arguments as your beliefs implies that you are not contradicting yourself. We further define a single argument A as being *acceptable with respect to a set of arguments* S , if for each set of arguments $T \subseteq A$, such that $T \triangleright A$, there is an argument B in T , such that S attacks B . In that case we also say that S *defends* A . A conflict-free set S , where all arguments in S are acceptable with respect to S , is called *admissible*. Thus, claiming an admissible set of arguments as your beliefs means that you can defend your beliefs against all possible counter arguments.

A very skeptical attitude towards arguments is captured by a semantics called the *grounded extension*. The grounded extension of an argumentation framework is the least fixpoint of the function $F : 2^A \rightarrow 2^A$, defined as

$$F(S) = \{A \in A : A \text{ is acceptable wrt. } S\}.$$

Thus the grounded extension consists of all the arguments that have no counter arguments, along with the arguments they defend, along with those defended by these arguments, and so on. A less skeptical semantics is that of a *preferred*

² Traditionally, attack relations have primarily been defined as subsets of $A \times A$. The generalized notion that we use here renders the language, argumentation system, and proofs to follow more elegant, but there are other, more general, reasons for preferring the definition used here. See [19] and [17] for dialectical and computational reasons.

extension, which is an admissible set that is maximal with respect to set inclusion. Finally, an admissible set S is said to be a *stable extension*, if it attacks all arguments in $A \setminus S$. Clearly, a stable extension is a preferred extension as well.

In general it is hard to compute a preferred extension [5], but Doutre and Mengin [6] present a method that enumerates preferred extensions for an abstract argumentation system as presented in [7]. Furthermore, Vreeswijk and Prakken [28] and Cayrol et al. [2] present methods for answering whether a specific argument is in at least one preferred extension, or if it is in all preferred extensions in the special case, where each preferred extension of the argumentation system is a stable one. In [18], we have adapted the technique of [6] to the problem of enumerating preferred extensions for argumentation systems, where sets of arguments attack other arguments. We present only the essentials of this technique here, and refer the interested reader to [18] for the details:

Given an argumentation system $\mathcal{A} \equiv (A, \triangleright)$, we define an \mathcal{A} -candidate as a triple $\mathcal{C} \equiv (I, O, U \equiv A \setminus (I \cup O))$ where

- $I \cap O = \emptyset$,
- every argument that is attacked by I is in O , and
- every argument A , for which there exists $S \subseteq I$ and $B \in I$, such that $S \cup A \triangleright B$, is in O .

Given an \mathcal{A} -candidate $\mathcal{C} \equiv (I, O, U)$ and an argument $A \in U$ the triples $\mathcal{C} - A \equiv (I_{-A}, O_{-A}, U_{-A} \equiv A \setminus (I_{-A} \cup O_{-A}))$ and $\mathcal{C} + A \equiv (I_{+A}, O_{+A}, U_{+A} \equiv A \setminus (I_{+A} \cup O_{+A}))$ are given by:

$$\begin{aligned} I_{-A} &\equiv I, & O_{-A} &\equiv O \cup A, \\ I_{+A} &\equiv I \cup A, & \text{and } O_{+A} &\equiv O \cup \Delta_{\mathcal{C}+A}, \end{aligned}$$

where

$$\begin{aligned} \Delta_{\mathcal{C}+A} &\equiv \{B \in U: \exists S \subseteq I, C \in I \\ &\quad \text{s.t. } S \cup A \triangleright B \vee S \cup B \triangleright A \\ &\quad \vee S \cup \{A, B\} \triangleright C \vee S \cup \{A, B\} \triangleright A\}. \end{aligned}$$

If A does not participate in a minimal attack on itself (which is the case for all arguments of the argumentation system we construct in this paper), then both $\mathcal{C} - A$ and $\mathcal{C} + A$ are \mathcal{A} -candidates themselves, and we can thus construct *candidate trees* where each node is an \mathcal{A} -candidate: Each \mathcal{A} -candidate \mathcal{C} has two children $\mathcal{C} - A$ and $\mathcal{C} + A$, for some arbitrary chosen A in U , except those candidates where $U = \emptyset$, which act as leaves in the tree. A candidate tree having candidate \mathcal{C} as root, is called a \mathcal{C} -tree.

It can be proven that if S^* is a preferred extension of \mathcal{A} , then there is a leaf $\mathcal{C} \equiv (I, O, \emptyset)$ of any $(\emptyset, \emptyset, A)$ -tree such that $I = S^*$. Conversely, for any leaf in a $(\emptyset, \emptyset, A)$ -tree, where I defends itself, I is admissible. It follows that, by constructing an arbitrary $(\emptyset, \emptyset, A)$ -tree, all preferred extensions can be enumerated.

Nielsen and Parsons [18] give a number of pruning rules for candidate trees. These can be used during construction to cut off branches that cannot contain leaves with preferred extensions.

3. Compromising on Bayesian networks

The problem we are addressing arises in a MAS containing a finite number of cooperating agents. We assume an ordering over the agents exists, so that we can refer to them by integers. Each agent i has a BN \mathcal{B}_i over a common set of domain variables V , which we assume to be implicit in the remainder of the text. For ease of exposition, we furthermore assume that an arbitrary but fixed total ordering \leadsto over the variables is known by all agents a priori. At some point agents 1 to k decide to pool their knowledge, as represented by \mathcal{B}_1 to \mathcal{B}_k , into a new BN \mathcal{B}_* , the *compromise* BN. Facilitating this task is the problem addressed here. We expect \mathcal{B}_1 to \mathcal{B}_k to be large but somewhat similar (as each describe relationships among the same variables), and therefore that having each agent communicate its entire model to each other agent is inefficient.

We shall focus solely on the graphical structure of \mathcal{B}_* (although we allow for non-structural aspects to act as guidelines in the construction of the structure). Therefore, the outcome of debate can be the full DAG of \mathcal{B}_* , its pattern, or its completed pattern. These should all be equivalent. However, as the next example shows, this is not entirely true.

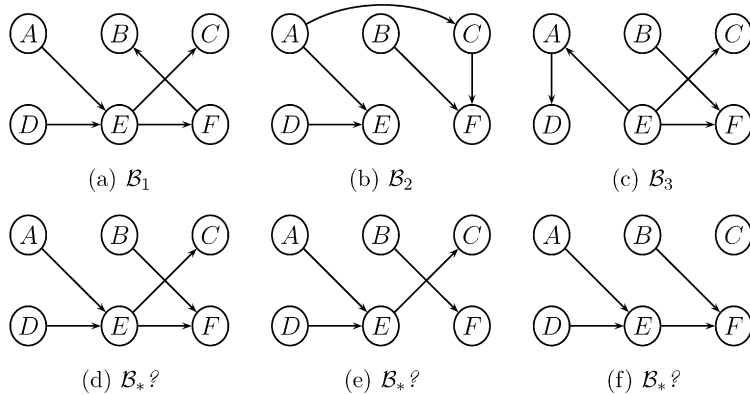


Fig. 6. Three BNs and three different compromise graphs.

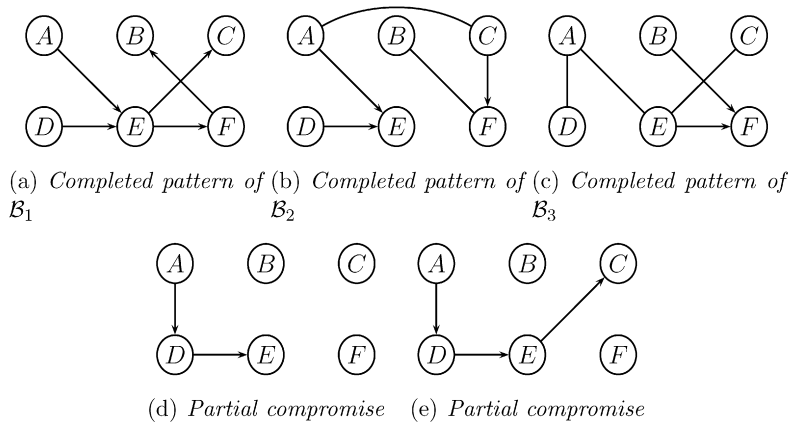


Fig. 7. Three completed patterns and two partial compromises.

Example 2 (*Agreeing on a compromise BN*). Consider three agents with BNs \mathcal{B}_1 , \mathcal{B}_2 , and \mathcal{B}_3 portrayed in Figs. 6(a) to 6(c). If the compromise is constructed by simple majority voting on the possible connections between each pair of nodes, the result is the structure in Fig. 6(d). If the same is done for the patterns of \mathcal{B}_1 to \mathcal{B}_3 instead, the BN structure in Fig. 6(e) is obtained. For the completed patterns the result is the one in Fig. 6(f).

As can be seen, evaluating each feature in isolation leads to vastly different results, even though the starting points are the same. This suggests that this simple approach is problematic, and should be ruled out, or that we should settle for one representation only. We chose the latter route, as we see no simple characterization of those combination methodologies that guarantee that equivalent results are obtained.

There seems to be little epistemological justification for choosing to compromise on a full graph, as most/all learning algorithms are unable to differentiate between two graphs belonging to the same equivalence class, and compromising on a DAG therefore seems to involve compromising on too much. Therefore, we choose to have agents compromise on a pattern or completed pattern. Unfortunately, we are unaware of any simple characterization of graphs that are patterns, so we choose completed patterns out of necessity. In summary, we shall aim for having the agents compromise on the completed pattern $\mathcal{G}_* \in \mathcal{C}^{cp}$ of \mathcal{B}_* .

To establish whether a graph is a good compromise for the agents, we need a measure for how well such graphs match each of \mathcal{B}_1 to \mathcal{B}_k . Furthermore, as we plan to build this compromise gradually, we wish for this measure to be relative to an already agreed *partial compromise*. The need for this should be clear from the following example.

Example 3 (*The need for partial compromises*). Consider again the setting from Example 2, and refer to the completed patterns of \mathcal{B}_1 to \mathcal{B}_3 in Figs. 7(a) to 7(c). It may be the case that agents have already agreed upon the connections in

the graph in Fig. 7(d), and now Agent 2 is asked to evaluate the compromise where a link is added between A and C . Obviously, this addition is consistent with Agent 2's beliefs, and thus must be valued highly. Consider then another situation, where the agents have agreed upon the connections in Fig. 7(e), and Agent 2 is again asked to evaluate the addition of a link between A and C . Now the addition is not as important, since the connection between A and C would have to be directed into C (to avoid a directed cycle), and thus create a v-structure (A, C, E) not present in \mathcal{B}_2 . Moreover, A and C are already connected (albeit through D and E) and the relationship between the two might be sufficiently represented by this connection.

In general, we cannot assume that a partially specified graph is suitable as representation of a partial compromise, as this might include agreements on what should *not* be part of the final compromise. To address this need and otherwise opting for an as general solution as possible, we shall take a partial compromise $\mathcal{P} \equiv (\mathbf{P}_+, \mathbf{P}_-)$ to be two sets of sentences in some language, where \mathbf{P}_+ describe aspects that should be true of the final compromise, and \mathbf{P}_- describe aspects that cannot be true.

For any three partial compromises \mathcal{P} , \mathcal{P}^a , and \mathcal{P}^b , where $\mathbf{P}_+ \subseteq \mathbf{P}_+^a$, $\mathbf{P}_- \subseteq \mathbf{P}_-^a$, $\mathbf{P}_+ \subseteq \mathbf{P}_+^b$ and $\mathbf{P}_- \subseteq \mathbf{P}_-^b$, we assume that each agent i can compute its *compromise scores* $s_i(\mathcal{P}, \mathcal{P}^a)$ and $s_i(\mathcal{P}, \mathcal{P}^b)$ such that $s_i(\mathcal{P}, \mathcal{P}^a) > s_i(\mathcal{P}, \mathcal{P}^b)$ iff \mathcal{P}^a describes \mathcal{B}_i better than \mathcal{P}^b , given that \mathcal{P} has already been accepted as being descriptive of \mathcal{B}_i . We will assume s_i to be additive, i.e. for any three partial compromises \mathcal{P}^0 , \mathcal{P}^1 , and \mathcal{P}^2 , where $\mathbf{P}_+^0 \subseteq \mathbf{P}_+^1 \subseteq \mathbf{P}_+^2$ and $\mathbf{P}_-^0 \subseteq \mathbf{P}_-^1 \subseteq \mathbf{P}_-^2$, it is the case that $s_i(\mathcal{P}^0, \mathcal{P}^2) = s_i(\mathcal{P}^0, \mathcal{P}^1) + s_i(\mathcal{P}^1, \mathcal{P}^2)$.

Example 4 (Compromise scoring functions). Consider the partial compromises $\mathcal{P}^a \equiv (\mathbf{P}_+^a, \mathbf{P}_-^a = \emptyset)$ and $\mathcal{P}^b \equiv (\mathbf{P}_+^b, \mathbf{P}_-^b = \emptyset)$ where the sentences in \mathbf{P}_+^a and \mathbf{P}_+^b represent the graphs in Figs. 7(d) and 7(e), respectively. A simple example of $s_2(\mathcal{P}^a, \mathcal{P}^b)$ could be the number of features described in $\mathbf{P}_+^b \setminus \mathbf{P}_+^a$, which are consistent with \mathcal{B}_2 , minus those that are not. Specifically, $s_2(\mathcal{P}^a, \mathcal{P}^b)$ would equal -1 .

A more complex score could weigh each of these described features according to the empirical evidence Agent 2 has in favor of or against them: We may have that Agent 2 is employing a sensor for measuring the E variable, and that this sensor is known to be of poor quality. Therefore, $s_2(\mathcal{P}^a, \mathcal{P}^b)$ would be equal to -0.3 , as an indication that the addition of the arc is against Agent 2's observations, but that those observations could easily be flawed.

Yet another score could take into account the ramifications of the sentences in $\mathbf{P}_+^b \setminus \mathbf{P}_+^a$. For instance, the addition of the arc from E to C might in itself not be much at odds with Agent 2's observations, but the addition forces it to either give up the idea of a link between A and C , or accept the existence of v-structure (A, C, E) , both of which it may have strong evidence for/against. Therefore, it has $s_2(\mathcal{P}^a, \mathcal{P}^b) = -10$. Of course, the challenge is for Agent 2 to actually be able to survey the impact of the sentences in $\mathbf{P}_+^b \setminus \mathbf{P}_+^a$.

Had \mathbf{P}_-^a not been empty, but rather included an agreement that there was not to be any connection between A and C , then the previous example of a compromise score would not have been so low, since the dreaded consequence of the addition of the arc from E to C would already be a consequence of \mathcal{P}^a . On the other hand, had this agreement been a part of $\mathbf{P}_-^b \setminus \mathbf{P}_-^a$, then the first two examples of scores would have been lower: -2 for the first, and something less than -0.3 , depending on Agent 2's evidence, for the second.

Notice, that with this open definition, we do not attempt to define what it means to be a “better description”, since we believe that this issue can be dependent on the actual setting in which the framework is to be used, as stated in Section 1.

In addition to the compromise score, we also assume that the agents know the *combination function* $c: \mathbb{R}^k \rightarrow \mathbb{R}$, indicating how much trust should be put into the individual agents' models. Differences in trust can be justified by differences in experiences and sensor accuracies. Formally, we define c as follows: Let \mathcal{P} , \mathcal{P}^a , and \mathcal{P}^b be partial compromises. c is the combination function for agents i to k , if

$$c(s_1(\mathcal{P}, \mathcal{P}^a), \dots, s_k(\mathcal{P}, \mathcal{P}^a)) > c(s_1(\mathcal{P}, \mathcal{P}^b), \dots, s_k(\mathcal{P}, \mathcal{P}^b)),$$

whenever \mathcal{P}^a is a better compromise than \mathcal{P}^b for the group of agents 1 to k , given that they have already agreed on \mathcal{P} . An obvious choice for c would be a linear combination of its inputs, and in fact we shall assume that this is the case in this text. We refer to $c(s_1(\mathcal{P}, \mathcal{P}^a), \dots, s_k(\mathcal{P}, \mathcal{P}^a))$ as the *joint compromise score* of \mathcal{P}^a given \mathcal{P} .

With this notation in place, we can restate the task more formally, as that of finding a partial compromise \mathcal{P} , which uniquely identifies some graph $\mathcal{G}_* \in \mathcal{C}^{cp}$, such that

$$c(s_1((\emptyset, \emptyset), \mathcal{P}), \dots, s_k((\emptyset, \emptyset), \mathcal{P})) \geq c(s_1((\emptyset, \emptyset), \mathcal{P}'), \dots, s_k((\emptyset, \emptyset), \mathcal{P}')),$$

for all other partial compromises \mathcal{P}' , which uniquely identifies a graph $\mathcal{G}' \in \mathcal{C}$.

As presented here, it is clear that the problem is not of a simple binary nature, as we are not trying to establish whether some proposition is true or not, and that we are furthermore dealing with a setting in which more than two agents may interact. Consequently, we cannot utilize the vast literature on dialectic proof theories directly. Rather, the problem we are trying to solve is a distributed maximization over a super exponential hypothesis space (\mathcal{C}). Furthermore, as the worth of (partial) compromises are specified in relation to already agreed upon compromises, the problem is of a highly dynamic nature.

Our solution to the problems is divided into three parts. First, we create a finite language with which graphs and some essential properties of these can be expressed; second and most importantly, we construct an argumentation system with which the agents can reason about consequences of committing to partial compromises; and thirdly, we create an agora in which the agents can reach compromise graphs in an anytime fashion.

4. Encoding graphs

For the agents to compromise on \mathcal{G}_* , a formal language L for expressing graphs and properties of graphs must be defined. For efficiency we aim to make this language as simple as possible, while ensuring that it is still sufficiently powerful to describe any graph and its membership status of \mathcal{C}^{cp} . By simple, we mean preferably finite and as small as possible.

First of all, we introduce a simple language L^g for encoding of graphs only:

Definition 5 (*Simple graph language*). The language L^g is the set such that $\text{Arc}(X, Y)$, $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, and $\text{NonAdjacent}(X, Y)$ are members of L^g iff X and Y ($X \rightsquigarrow Y$) are distinct variables.

A *graph knowledge base* we define to be a set $\Sigma^g \subseteq L^g$. Further:

Definition 6 (*Consistent and closed graph knowledgebases*). Given a graph knowledge base Σ^g , if it holds that for all pairs of variables X and Y , where $X \rightsquigarrow Y$, a maximum of one of $\text{Arc}(X, Y)$, $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, and $\text{NonAdjacent}(X, Y)$ is in Σ^g , then we call Σ^g a *consistent graph knowledgebase (CGK)*.

Furthermore, if it holds that for any two variables X and Y , where $X \rightsquigarrow Y$, exactly one of $\text{Arc}(X, Y)$, $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, and $\text{NonAdjacent}(X, Y)$ is in Σ^g , then we call Σ^g a *closed graph knowledgebase (CLGK)*.

The *graph encoded by* a CGK Σ^g is the graph $\mathcal{G}[\Sigma^g]$ resulting from starting with the graph with no edges, and then for any two nodes X and Y ($X \rightsquigarrow Y$) adding an arc from X to Y if $\text{Arc}(X, Y)$ is in Σ^g , an arc from Y to X if $\text{Arc}(Y, X)$ is in Σ^g , or an undirected edge if $\text{Link}(X, Y)$ is in Σ^g . It is easy to see that graph encoded by a CGK is well-defined. Furthermore, given a graph $\mathcal{G} \equiv (V, E)$ there exists a unique CLGK $\Sigma^g[\mathcal{G}]$, for which \mathcal{G} is the encoded graph. Clearly, $\Sigma^g[\mathcal{G}]$ can be constructed in time $O(|V|^2)$.

Example 7 (*Graph encoding*). Consider the graph knowledgebase Σ_e^g comprised of $\text{Arc}(A, D)$, $\text{Arc}(D, E)$, and $\text{Arc}(E, C)$. Σ_e^g is a CGK, albeit not a closed one, as the connections among several pairs of variables, such as A and E are not specified. Assuming that L^g is built from variables $\{A, B, C, D, E, F\}$, $\mathcal{G}[\Sigma_e^g]$ corresponds to the graph in Fig. 7(e). If furthermore all of

$\text{NonAdjacent}(A, B)$, $\text{NonAdjacent}(A, C)$, $\text{NonAdjacent}(A, E)$,
 $\text{NonAdjacent}(A, F)$, $\text{NonAdjacent}(B, C)$, $\text{NonAdjacent}(B, D)$,
 $\text{NonAdjacent}(B, E)$, $\text{NonAdjacent}(B, F)$, $\text{NonAdjacent}(C, D)$,
 $\text{NonAdjacent}(C, F)$, $\text{NonAdjacent}(D, F)$, and $\text{NonAdjacent}(E, F)$

were added to Σ_e^g , Σ_e^g would be the encoded version of the graph in Fig. 7(e), and would therefore be a CLGK. If any further sentence was added to Σ_e^g , the set would cease to be a CGK.

We thus have that any graph can be efficiently encoded by a unique CLGK, and Definition 6 allows us to distinguish the graph knowledge bases, which can be interpreted as graphs, from those that cannot. Next, we extend L^g into a language powerful enough for building a reasoning engine about graphs and their membership status of \mathcal{C}^{cp} on top:

Definition 8 (*Graph language*). The graph language L is the set consisting of all sentences in L^g and

- $\text{ArcNotAllowed}(X, Y)$,
- $\text{DirectedPath}(X, Y)$,
- $\text{UndirectedPath}(X, Y)$,
- $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$,
- $\neg \text{DirectedPath}(X, Y)$,
- $\neg \text{UndirectedPath}(X, Y)$, and
- $\neg \text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$,

for any choice of distinct variables X, Y, Z , and W ³ ($Z \rightsquigarrow W$). Sentences like $\text{DirectedPath}(X, Y)$, $\text{UndirectedPath}(X, Y)$, and $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ will be referred to as *path sentences*, and sentences like $\neg \text{DirectedPath}(X, Y)$, $\neg \text{UndirectedPath}(X, Y)$, and $\neg \text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ will be referred to as *negative path sentences*.

Intuitively, the sentences just introduced are supposed to be used as descriptors of attributes of the graphs encoded by CGKs: $\text{ArcNotAllowed}(X, Y)$ states that an arc from X to Y would not be strongly protected; $\text{DirectedPath}(X, Y)$ states that there is a directed path from X to Y ; $\text{UndirectedPath}(X, Y)$ states that there is an undirected path between X and Y ; $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ states that there is an undirected path not comprising Z nor W between X and Y ; $\neg \text{DirectedPath}(X, Y)$ states that there is no directed path from X and Y ; $\neg \text{UndirectedPath}(X, Y)$ states that there is no undirected path between X and Y ; and $\neg \text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ states that there is no undirected path between X and Y , or that any such path necessarily contains either Z or W .

As L^g is a subset of L , it follows that a graph knowledge base is a set of sentences in L as well, and in particular that Definition 6 still makes sense. Given a set Σ of sentences of L , we denote by Σ^g the set $\Sigma \cap L^g$.

5. Graph argumentation system

Building on the language L introduced above, we define an argumentation system for distinguishing completed patterns that could be compromises for the agents. The system that we construct enjoys the properties that a graph is a member of \mathcal{C} iff there is a preferred extension of the system, such that the extension encodes this graph.

Definition 9 (*Graph argumentation system*). The graph argumentation system \mathcal{A}^g is the tuple $(L, \triangleright^g \subseteq (2^L \times L))$, where \triangleright^g is defined as follows ($[A, B]$ is short-hand for any one of (A, B) and (B, A)):

- (1) $\text{Arc}(X, Y) \triangleright^g \text{Arc}(Y, X)$
- (2) $\text{Arc}(X, Y) \triangleright^g \text{Link}[X, Y]$
- (3) $\text{Arc}(X, Y) \triangleright^g \text{NonAdjacent}[X, Y]$
- (4) $\text{Link}(X, Y) \triangleright^g \text{Arc}[X, Y]$
- (5) $\text{Link}(X, Y) \triangleright^g \text{NonAdjacent}[X, Y]$
- (6) $\text{NonAdjacent}(X, Y) \triangleright^g \text{Arc}[X, Y]$
- (7) $\text{NonAdjacent}(X, Y) \triangleright^g \text{Link}[X, Y]$
- (8) $\neg \text{DirectedPath}(X, Y) \triangleright^g \text{DirectedPath}(X, Y)$
- (9) $\neg \text{UndirectedPath}(X, Y) \triangleright^g \text{UndirectedPath}(X, Y)$

³ Throughout the text we assume that the implicit set of variables has at least five members. This assumption can be lifted, albeit with a more complex notation to follow.

- (10) $\neg \text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W) \triangleright^g \text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$
- (11) $\text{Arc}(X, Y) \triangleright^g \neg \text{DirectedPath}(X, Y)$
- (12) $\text{Link}(X, Y) \triangleright^g \neg \text{UndirectedPath}[X, Y]$
- (13) $\text{Link}(X, Y) \triangleright^g \neg \text{UndirectedPath}[X, Y] \text{ Excluding}(Z, W)$
- (14) $\{\text{DirectedPath}(X, Y), \text{DirectedPath}(Y, Z)\} \triangleright^g \neg \text{DirectedPath}(X, Z)$
- (15) $\{\text{DirectedPath}(X, Y), \text{UndirectedPath}[Y, Z]\} \triangleright^g \neg \text{DirectedPath}(X, Z)$
- (16) $\{\text{UndirectedPath}[X, Y], \text{DirectedPath}(Y, Z)\} \triangleright^g \neg \text{DirectedPath}(X, Z)$
- (17) $\{\text{UndirectedPath}[X, Y], \text{UndirectedPath}[Y, Z]\} \triangleright^g \neg \text{UndirectedPath}[X, Z]$
- (18) $\{\text{UndirectedPath}[X, Y] \text{ Excluding}(Z, W), \text{UndirectedPath}[Y, U] \text{ Excluding}(Z, W)\} \triangleright^g \neg \text{UndirectedPath}[X, U] \text{ Excluding}(Z, W)$
- (19) $\text{DirectedPath}(X, Y) \triangleright^g \text{Arc}(Y, X)$
- (20) $\text{DirectedPath}(X, Y) \triangleright^g \text{Link}[X, Y]$
- (21) $\text{UndirectedPath}[X, Y] \triangleright^g \text{Arc}(X, Y)$
- (22) $\{\text{UndirectedPath}[X, Y] \text{ Excluding}(W, Z), \text{Link}[X, W], \text{Link}[Y, Z], \text{NonAdjacent}[X, Z], \text{NonAdjacent}[Y, W]\} \triangleright^g \text{Link}[W, Z]$
- (23) $\{\text{Arc}(X, Y), \text{NonAdjacent}[X, Z]\} \triangleright^g \text{Link}[Y, Z]$
- (24) $\text{ArcNotAllowed}(X, Y) \triangleright^g \text{Arc}(X, Y)$
- (25) $\{\text{Arc}(Z, X), \text{NonAdjacent}[Z, Y]\} \triangleright^g \text{ArcNotAllowed}(X, Y)$
- (26) $\{\text{Arc}(Z, Y), \text{NonAdjacent}[Z, X]\} \triangleright^g \text{ArcNotAllowed}(X, Y)$
- (27) $\{\text{Arc}(X, Z), \text{Arc}(Z, Y)\} \triangleright^g \text{ArcNotAllowed}(X, Y)$
- (28) $\{\text{Link}[X, Z], \text{Arc}(Z, Y), \text{Link}[X, W], \text{Arc}(W, Y), \text{NonAdjacent}[Z, W]\} \triangleright^g \text{ArcNotAllowed}(X, Y)$

for all choices of distinct variables X, Y, Z, W , and U where the sentences obtained are in L .

Loosely speaking, if Σ is a conflict free set wrt. \triangleright^g , then items 1–7 ensure that Σ^g is a CGK; items 8–18 make sure that the path and negative path sentences in $\Sigma \setminus \Sigma^g$ are correct wrt. the graph that is encoded by Σ^g ; items 19–21 ensure that Σ^g encodes a chain graph; item 22 ensures that the graph encoded by Σ^g has decomposable chain components; item 23 ensures that graph also respects item 3 of Theorem 1; and items 24–28 guarantee that all arcs in the graph are strongly protected.

Example 10 (*Argumentative reasoning about graphs*). We return to the CGK Σ_e^g presented in Example 7 joined with sentences $\text{DirectedPath}(A, D), \text{DirectedPath}(D, E), \text{DirectedPath}(E, C), \text{DirectedPath}(D, C)$, and $\text{DirectedPath}(A, C)$ to constitute the set of sentences Σ_e . Clearly, Σ_e is conflict-free. From item 19 we can conclude that $\text{Arc}(C, A)$ cannot be added to Σ_e without destroying its property of being conflict-free wrt. \mathcal{A}^g . Item 20 yields the same result for $\text{Link}(A, C)$, and the only options for putting a description of the connection between A and C into Σ_e are thus $\text{NonAdjacent}(A, C)$ and $\text{Arc}(A, C)$. The argumentation system thus allows for reasoning about such consequences that were treated informally in Examples 3 and 4.

In what follows we provide a formal treatment of the above outlined intuitions. More specifically we prove two important results, namely that any preferred extension of \mathcal{A}^g is a member of \mathfrak{C}^{cp} , and that any member of \mathfrak{C}^{cp} has a corresponding stable extension of \mathcal{A}^g . These results are important since they guarantee that, if agents seek compromises under the restrictions specified by \mathcal{A}^g , they can be sure that their result is a completed pattern and that they are not restricted from agreeing on any model a priori by the relations of \mathcal{A}^g .

Lemma 11. *Let Σ be conflict free wrt. \mathcal{A}^g . Then Σ^g is a CGK.*

Proof. Obvious, given items 1–7. \square

Lemma 12. *Let Σ be a preferred extension of \mathcal{A}^g . Then Σ^g is a CLGK.*

Proof. We need to prove that if X and Y ($X \rightsquigarrow Y$) are two nodes that are non-adjacent in $\mathcal{G}[\Sigma^g]$ then $\text{NonAdjacent}(X, Y)$ is in Σ . We prove this by contradiction: We assume that $\text{NonAdjacent}(X, Y)$ is not in Σ and then prove that $\Sigma^* = \Sigma \cup \text{NonAdjacent}(X, Y)$ is an admissible set, which contradicts that Σ is a preferred extension.

First we prove that Σ^* is conflict-free: Since X and Y are non-adjacent in $\mathcal{G}[\Sigma^g]$ it follows that none of $\text{Arc}(X, Y)$, $\text{Arc}(Y, X)$, and $\text{Link}(X, Y)$ are in Σ , and therefore that Σ cannot attack $\text{NonAdjacent}(X, Y)$. Thus, if there is a set $S \subseteq \Sigma^*$ and an argument $A \in \Sigma^*$, such that $S \triangleright^g A$, then A cannot be $\text{NonAdjacent}(X, Y)$ and hence must be in Σ . As Σ is admissible, it follows that Σ attacks at least one argument B in S , and as Σ does not attack $\text{NonAdjacent}(X, Y)$, B must be a member of $S \setminus \text{NonAdjacent}(X, Y)$. But that means that Σ attacks Σ , meaning that it is not a conflict-free set, contradicting the assumption that Σ is a preferred extension. So we conclude that Σ^* must be conflict-free.

It is easy to see that Σ^* attacks all sets of arguments that attack Σ^* : By assumption, Σ attacks all sets of arguments that attack an argument in Σ , and by items 1–7, $\text{NonAdjacent}(X, Y)$ attacks all sets of arguments that attack itself. Consequently, $\Sigma^* = \Sigma \cup \text{NonAdjacent}(X, Y)$ attacks all arguments that attack some argument in itself, and the lemma follows. \square

From [19] we have the following equivalent of the “Fundamental Lemma” of [7]:

Lemma 13. *Let Σ be a preferred extension, and let A an argument defended by Σ . Then A is in Σ .*

Lemma 14. *Let Σ be a preferred extension of \mathcal{A}^g and X, Y, Z , and W ($Z \rightsquigarrow W$) be variables. Then*

- (1) *if there is a directed path from X to Y in $\mathcal{G}[\Sigma^g]$, then $\text{DirectedPath}(X, Y)$ is in Σ ,*
- (2) *if there is an undirected path from X to Y in $\mathcal{G}[\Sigma^g]$, then $\text{UndirectedPath}(X, Y)$ is in Σ , and*
- (3) *if there is an undirected path from X to Y in $\mathcal{G}[\Sigma^g]$ not comprising any of the variables Z and W , then $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ is in Σ .*

Proof. Induction on the length l of the path between X and Y . We first consider the base case where l is 1. If the path is a directed path, this means that there is an arc from X to Y and consequently that $\text{Arc}(X, Y)$ is in Σ . By item 11 this means that Σ attacks $\neg \text{DirectedPath}(X, Y)$, which is the only argument that attacks $\text{DirectedPath}(X, Y)$. Lemma 13 then guarantees that $\text{DirectedPath}(X, Y)$ is in Σ . Similar arguments establish the base cases for 2 and 3.

For the induction step, assume that the result is valid for all paths of length $l \leq n$ and consider a path $\pi = (X = X_1, U = X_2, \dots, X_{n+2} = Y)$ of length $n + 1$. If π is an undirected path, then we know from the induction hypothesis that both $\text{UndirectedPath}(X, U)$ and $\text{UndirectedPath}(U, Y)$ must be in Σ . The two sentences collectively attack $\neg \text{UndirectedPath}(X, Y)$, which is the only argument attacking $\text{UndirectedPath}(X, Y)$. Since Σ is a preferred extension, Lemma 13 then gives that $\text{UndirectedPath}(X, Y)$ is a member of Σ . The proof for $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ is the same.

Assume next that π is a directed path. Then either there is an arc from X to U or a link between them. We consider the two cases in turn.

In the first case, an argument similar to the one used for the base case gives us that $\text{DirectedPath}(X, U)$ must be in Σ . Next, if $(U = X_2, \dots, X_{n+2} = Y)$ is a directed path, then induction hypothesis gives us that $\text{DirectedPath}(U, Y)$ is in Σ . If $(U = X_2, \dots, X_{n+1} = Y)$ is undirected, we get that $\text{UndirectedPath}(U, Y)$ is in Σ . Either way, the two sentences obtained collectively attack $\neg \text{DirectedPath}(X, Y)$, which is the only argument attacking $\text{DirectedPath}(X, Y)$. Since Σ is a preferred extension, Lemma 13 then states that $\text{DirectedPath}(X, Y)$ must be a member of Σ .

If there is a link between X and U , there must be a directed path from U to Y , for π to be a directed path. As before the presence of $\text{Link}(X, U)$ or $\text{Link}(U, X)$ yields that $\text{UndirectedPath}(X, U)$ is in Σ , and the induction hypothesis that $\text{DirectedPath}(U, Y)$ is in Σ . Together the two attack $\neg \text{DirectedPath}(X, Y)$, which is the only argument that attacks $\text{DirectedPath}(X, Y)$. Once again Lemma 13 guarantees that $\text{DirectedPath}(X, Y)$ is in Σ . \square

The converse of Lemma 14 is not true, which can be seen from the following example:

Example 15. Let Σ^g consist of $\text{NonAdjacent}(A, B)$, $\text{NonAdjacent}(B, C)$, and $\text{NonAdjacent}(A, C)$, thus encoding the empty graph over the set of variables $V_e = \{A, B, C\}$. Furthermore, let $\Sigma \setminus \Sigma^g$ be

$$\begin{aligned} & \{\text{ArcNotAllowed}(X, Y): X, Y \in V_e\} \\ & \cup \{\text{UndirectedPath}(X, Y): X, Y \in V_e\} \\ & \cup \{\text{DirectedPath}(X, Y): X, Y \in V_e\}. \end{aligned}$$

Then Σ is conflict free and attacks all other sentences in L , meaning that it is stable and hence a preferred extension, yet clearly $\mathcal{G}[\Sigma^g]$ does not reflect the meaning of the path sentences in Σ .

Theorem 16. *Let Σ be a preferred extension of \mathcal{A}^g . Then $\mathcal{G}[\Sigma^g]$ is in \mathcal{C}^{cp} .*

Proof. Assume not. This means that $\mathcal{G}[\Sigma^g]$ fails to satisfy one of the bullets of Theorem 1. We prove that this cannot be the case, one bullet at the time:

Assume $\mathcal{G}[\Sigma^g]$ does not satisfy item 1, namely that it is not a chain graph. Then there is a directed cycle (X, \dots, Y, X) in $\mathcal{G}[\Sigma^g]$ containing at least one arc. Without loss of generality, assume this arc to be going from Y to X , and consequently that $\text{Arc}(Y, X)$ is in Σ . The rest of the cycle then constitutes either a directed path from X to Y , or an undirected path between them. Either way, Lemma 14 allows us to conclude that one of $\text{DirectedPath}(X, Y)$ and $\text{UndirectedPath}[X, Y]$ must be in Σ . But that is impossible, since either one of them attacks $\text{Arc}(Y, X)$ and Σ is conflict-free.

Assume then that $\mathcal{G}[\Sigma^g]$ does not satisfy item 2, namely that it contains a chain component that is not decomposable. This implies that there is a chordless cycle in $\mathcal{G}[\Sigma^g]$. Let this cycle be (X, Z, \dots, W, Y, X) , which implies that $\text{Link}[X, Z]$, $\text{Link}[W, Y]$, and $\text{Link}[Y, X]$ are in Σ . Furthermore, as the cycle is chordless, both X and W and Y and Z must be pairs of non-adjacent nodes in $\mathcal{G}[\Sigma^g]$. By Lemma 12 $\text{NonAdjacent}[X, W]$ and $\text{NonAdjacent}[Y, Z]$ must then be in Σ . Additionally, the path (Z, \dots, W) and Lemma 14 collectively allow us to conclude that $\text{UndirectedPath}[Z, W] \text{Excluding}(X, Y)$ must be in Σ as well. However, $\text{UndirectedPath}[Z, W] \text{Excluding}(X, Y)$ in conjunction with $\text{Link}[Z, X]$, $\text{Link}[W, Y]$, $\text{NonAdjacent}[X, W]$, and $\text{NonAdjacent}[Y, Z]$ attacks $\text{Link}[X, Y]$, which is impossible since Σ is conflict-free.

Next, assume that $\mathcal{G}[\Sigma^g]$ does not satisfy item 3, which means that there are three variables X , Y , and Z , such that X is a parent of Y , X and Z are not adjacent, and Z is a neighbor of Y . According to Lemma 12, Σ would then contain all of $\text{Arc}(X, Y)$, $\text{NonAdjacent}[X, Z]$, and $\text{Link}[Y, Z]$. But the first two attack the last one, and since Σ is conflict-free, this is impossible.

Finally, assume that $\mathcal{G}[\Sigma^g]$ does not satisfy item 4, which means that there is some arc from a node X to a node Y that is not strongly protected in $\mathcal{G}[\Sigma^g]$. Thus $\text{Arc}(X, Y)$ is in Σ , and as $\text{Arc}(X, Y)$ is attacked by $\text{ArcNotAllowed}(X, Y)$ and Σ is an admissible set, it follows that Σ attacks $\text{ArcNotAllowed}(X, Y)$. The only ways Σ can attack $\text{ArcNotAllowed}(X, Y)$ are by containing either

- $\text{Arc}(Z, X)$ and $\text{NonAdjacent}[Z, Y]$ for some variable Z ,
- $\text{Arc}(Z, Y)$ and $\text{NonAdjacent}[Z, X]$ for some variable Z ,
- $\text{Arc}(X, Z)$ and $\text{Arc}(Z, Y)$ for some variable Z , or
- $\text{Link}[X, Z]$, $\text{Link}[X, W]$, $\text{Arc}(Z, Y)$, $\text{Arc}(W, Y)$, and $\text{NonAdjacent}(Z, W)$ for some variables Z and W .

We deal only with the first case as the others are similar. As Σ contains $\text{Arc}(Z, X)$ and $\text{NonAdjacent}[Z, Y]$ and Σ^g is a CGK, it follows that there is an arc from Z to X in $\mathcal{G}[\Sigma^g]$ and that Z and Y are not connected by a link nor an arc. But then the arc from X to Y is protected, which yields a contradiction with the assumption. As the other three cases yield similar contradictions, the assumption must be false, and the theorem follows. \square

Theorem 17. *If \mathcal{G} is in \mathcal{C}^{cp} , then there is a stable extension Σ of \mathcal{A}^g , such that $\mathcal{G}[\Sigma^g] = \mathcal{G}$.*

Proof. We prove this by construction of a stable extension Σ such that $\mathcal{G}[\Sigma^g] = \mathcal{G}$. First let Σ^g be $\Sigma^g[\mathcal{G}]$, which ensures that $\mathcal{G}[\Sigma^g] = \mathcal{G}$, no matter what sentences are in $\Sigma \setminus \Sigma^g$.

Next, for any four variables X , Y , Z , and W , let $\Sigma \setminus \Sigma^g$ contain

- $\text{ArcNotAllowed}(X, Y)$ iff there is no arc between X and Y and such an arc would not be strongly protected,
- $\text{DirectedPath}(X, Y)$ iff there is a directed path from X to Y , and $\neg \text{DirectedPath}(X, Y)$ otherwise,
- $\text{UndirectedPath}(X, Y)$ iff there is an undirected path between X and Y , and $\neg \text{UndirectedPath}(X, Y)$ otherwise, and
- $\text{UndirectedPath}(X, Y) \text{ Excluding } (Z, W)$ iff there is an undirected path not comprising Z nor W between X and Y , and $\neg \text{UndirectedPath}(X, Y) \text{ Excluding } (Z, W)$ otherwise.

Then Σ is the stable extension we are seeking. To prove it, we show (i) that Σ is conflict-free and (ii) that Σ attacks all sentences in $L \setminus \Sigma$.

(i) We consider each bullet of the definition of \triangleright^g in Definition 9.

Items 1 to 7: As Σ^g was defined to be $\Sigma^g[\mathcal{G}]$, which is a CLGK, there can be no two variables X and Y ($X \rightsquigarrow Y$), for which more than one of $\text{Arc}(X, Y)$, $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, and $\text{Link}(Y, X)$ are in Σ . Consequently, these bullets do not give rise to any conflicts.

Items 8 to 13: These bullets cannot give rise to any conflict. The first because by definition of Σ , for each pair of variables X and Y , only one of $\text{DirectedPath}(X, Y)$ and $\neg \text{DirectedPath}(X, Y)$ can be in Σ . The same reasoning prevents the next two bullets from giving rise to a conflict. Item 9 cannot give rise to a conflict as the presence of $\text{Arc}(X, Y)$ in Σ means that there is an arc from X to Y in \mathcal{G} and thus a directed path from X to Y , which by definition of Σ means that it does not contain $\neg \text{DirectedPath}(X, Y)$. The same reasoning goes for the two next bullets.

Items 14 to 18 cannot give rise to conflicts. If item 14 was the reason for a conflict, then both $\text{DirectedPath}(X, Y)$ and $\text{DirectedPath}(Y, Z)$ would have to be in Σ , which by definition of Σ would mean that there is a directed path (X, \dots, Y) from X to Y and a directed path (Y, \dots, Z) from Y to Z in \mathcal{G} . But then there is also a directed path (X, \dots, Y, \dots, Z) from X to Z meaning that $\neg \text{DirectedPath}(X, Z)$ by definition is not in Σ . Similar arguments show that the other bullets cannot give rise to conflicts.

Item 19 cannot give rise to a conflict as that would mean that both $\text{DirectedPath}(X, Y)$ and $\text{Arc}(Y, X)$ would be in Σ , and consequently that there would be a directed path from X to Y and an arc from Y to X in \mathcal{G} , which would mean that \mathcal{G} is not a chain graph, and by Theorem 1 not a completed pattern. The same goes for items 20 to 21, and a similar argument, substituting the chain graph requirement of Theorem 1 with the decomposability one, yields that item 22 cannot be cause of conflicts either.

Item 23 cannot give rise to a conflict, as if it did that would mean that all of $\text{Arc}(X, Y)$, $\text{NonAdjacent}[X, Z]$, and $\text{Link}[Y, Z]$ would be in Σ^g , which in turn would mean that X would be a parent of Y in \mathcal{G} , that X and Z would not be adjacent, and that Y and Z would be connected with an undirected link. This would violate item 3 of Theorem 1, meaning that \mathcal{G} is not a completed pattern, which is a contradiction.

Item 24 can clearly not give rise to any conflicts, as $\text{ArcNotAllowed}(X, Y)$ are only in Σ for nodes not connected by an arc, ruling out that $\text{Arc}(X, Y)$ is also in Σ , which would be needed for a conflict.

The arguments as to why items 25 to 28 cannot give rise to conflicts are similar so only the argument relating to item 25 is treated here. If this bullet caused a conflict, it means that $\text{Arc}(Z, X)$, $\text{NonAdjacent}[Z, Y]$, and $\text{ArcNotAllowed}(X, Y)$ are all in Σ , implying that Z is a parent of X and not adjacent to Y in \mathcal{G} , and that there is no arc between X and Y , and had there been such an arc, it would not be strongly protected. But the definition of strongly protected clearly states that an arc from X to Y with Z being a parent of X and not adjacent to Y is strongly protected, yielding a contradiction.

(ii) Next, we prove that Σ attacks every sentence not in Σ , and thus that it is stable and therefore a preferred extension. We consider each argument A of L in turn, and show that it is either a member of Σ or attacked by it.

First, let A be $\text{Arc}(X, Y)$ for some variables X and Y ($X \rightsquigarrow Y$). If $\text{Arc}(X, Y)$ is not in Σ^g then, by definition of a CLGK either $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, or $\text{NonAdjacent}(X, Y)$ must be in Σ^g and Σ thus attacks $\text{Arc}(X, Y)$ by one of items 1 to 7. A similar argument applies if A is $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, or $\text{NonAdjacent}(X, Y)$.

Assume then that A is $\text{ArcNotAllowed}(X, Y)$ for some pair of variables X and Y , and that Σ does not attack $\text{ArcNotAllowed}(X, Y)$. That means that none of the left hand sides of items 25 to 28 are satisfied by members of Σ . This in turn means that the graph encoded by Σ^g , viz. \mathcal{G} , does not give strong protection to an arc from X to Y . As \mathcal{G} is a completed pattern, it follows from Theorem 1 that there cannot be an arc from X to Y . However, by definition of Σ this means that $\text{ArcNotAllowed}(X, Y)$ must be in Σ .

Then assume that A is $\text{DirectedPath}(X, Y)$ and that it is not attacked by Σ . This means that $\neg\text{DirectedPath}(X, Y)$ is not in Σ , and by definition of Σ , $\text{DirectedPath}(X, Y)$ must therefore be in Σ . Similar reasoning establishes that both $\text{UndirectedPath}(X, Y)$ and $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ are in Σ iff they are not attacked by Σ .

Finally, assume that A is $\neg\text{UndirectedPath}(X, Y)$, and that it is not in Σ . By definition of Σ this means that there is an undirected path (X, \dots, Y) between X and Y in \mathcal{G} . If the length of (X, \dots, Y) is 1, then there is an undirected link between the two variables X and Y , and consequently $\text{Link}[X, Y]$ is in Σ^g and Σ therefore attacks $\neg\text{UndirectedPath}(X, Y)$ by Item 12. If the length of (X, \dots, Y) is more than 1, there must be at least one variable Z on this path such that both (X, \dots, Z) and (Z, \dots, Y) are undirected paths in \mathcal{G} . It follows by the definition of Σ that both $\text{UndirectedPath}(X, Z)$ and $\text{UndirectedPath}(Z, Y)$ are in Σ , and Σ then attacks $\neg\text{UndirectedPath}(X, Y)$ by item 17. Similar (but more elaborate) arguments establish the result for arguments A of the form $\neg\text{DirectedPath}(X, Y)$ and $\neg\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$. \square

Given these results it is thus clear that the result Σ of a debate, if undertaken respecting \mathcal{A}^g , is a completed pattern if it is a preferred extension of \mathcal{A}^g , and that Σ can represent any completed pattern. However, checking whether a set of arguments constitute a preferred extension is complex. It involves checks for both admissibility and maximality. We therefore end this section with a result that yields a computationally efficient way of testing whether a set of arguments of \mathcal{A}^g is a preferred extension.

Lemma 18. *Let Σ be a preferred extension of \mathcal{A}^g , and A a path sentence not in Σ and not attacked by it. Then the negative path sentence corresponding to A is not in Σ nor attacked by it.*

Proof. We show only the case where the path sentence is $\text{UndirectedPath}(X, Y)$ for two variables X and Y , as the others cases are similar. Since $\text{UndirectedPath}(X, Y)$ is not in Σ and not attacked by Σ , it follows from item 9, and the fact that this bullet represents the only attack on $\text{UndirectedPath}(X, Y)$, that $\neg\text{UndirectedPath}(X, Y)$ cannot be in Σ either. Furthermore, as Σ is a preferred extension and does not attack $\text{UndirectedPath}(X, Y)$, it follows that Σ cannot attack $\neg\text{UndirectedPath}(X, Y)$ either, as that would imply that $\text{UndirectedPath}(X, Y)$ was defended by Σ , which according to Lemma 13 would mean that it should have been a member of Σ . \square

Theorem 19. *Let Σ be a preferred extension of \mathcal{A}^g . Then Σ is a stable extension.*

Proof. We need to show that Σ attacks each argument A in $\mathcal{A} \setminus \Sigma$. We do so by considering each possible A :

Assume first that A is one of $\text{Arc}(X, Y)$, $\text{Arc}(Y, X)$, $\text{Link}(X, Y)$, or $\text{NonAdjacent}(X, Y)$, for some variables X and Y . As Σ is a preferred extension, Lemma 12 guarantees that one of the other three arguments is in Σ , and items 1 to 7 then ensures that Σ attacks A .

Assume then that A is $\text{ArcNotAllowed}(X, Y)$, for some X and Y , and that it is not attacked by Σ . This means that for each of the left-hand sides of items 25 to 28 there is at least one element that is not in Σ . Furthermore, as Σ does not defend $\text{ArcNotAllowed}(X, Y)$, there is at least one of these bullets, whose left-hand side is not attacked by Σ . Assume that this bullet is item 25. That means that either $\text{Arc}(Z, X)$ or $\text{NonAdjacent}[Z, Y]$ or both are neither in Σ nor attacked by it. But this is impossible, as was proved above. A similar contradiction arise for the remaining three bullets. It follows that $\text{ArcNotAllowed}(X, Y)$ must be attacked by Σ .

Now, assume that A is $\neg\text{UndirectedPath}(X, Y)$, and that Σ does not attack $\neg\text{UndirectedPath}(X, Y)$. As Σ does not attack $\neg\text{UndirectedPath}(X, Y)$, it follows from item 17 that for all variables Z , one of $\text{UndirectedPath}[X, Z]$ and $\text{UndirectedPath}[Z, Y]$ cannot be in Σ either. Furthermore, item 12 implies that $\text{Link}[X, Y]$ cannot be in Σ , and since it is a preferred extension, and by Lemma 12 thus a CLGK, it must attack $\text{Link}[X, Y]$. Since by Lemma 13 Σ fails to defend $\neg\text{UndirectedPath}(X, Y)$, then it must fail to attack both of $\text{UndirectedPath}[X, Z]$ and $\text{UndirectedPath}[Z, Y]$ for some variable Z . Without loss of generality, assume that $\text{UndirectedPath}[Z, Y]$ is the path sentence which is both outside of Σ and not attacked by it. Lemma 18 then guarantees that $\neg\text{UndirectedPath}[Z, Y]$ must also be outside of Σ and not attacked by Σ . But $\neg\text{UndirectedPath}[Z, Y]$ has the exact same form as the A we started out exploring, and we can therefore apply the same argument again, obtaining yet another negative path sentence that must be outside Σ and not be attacked

by it. This process will never end, and since L is a finite language, it thus follows that there is a set N of negative path sentences, which are all outside Σ and not attacked by Σ solely because of other negative path sentences that are in N . From the definition of N it is obvious that the set P of path sentences, corresponding to the negative ones in N , attacks each of the sentences in N , and furthermore that N contains all the sentences that attack P . As Σ defends itself, it follows that the set $\Sigma^* \equiv \Sigma \cup P$ defends itself against all attacks. If we can also show that Σ^* is conflict-free, then it contradicts that Σ is a preferred extension, and hence implies that the original claim, that $\neg \text{UndirectedPath}(X, Y)$ is not attacked by Σ , is false. To show that Σ^* is conflict-free, we assume otherwise, and let $S \subseteq \Sigma^*$ attack some B in Σ^* . Assume first that B is in Σ . As Σ defends itself, it must attack some argument in S . But this is impossible since Σ is conflict-free and by construction P is not attacked by Σ . Thus, B must be in P . But this is also impossible, because the only arguments that attack P are the ones in N , which by definition is not part of Σ^* . It follows that Σ^* is conflict-free and thus that $\neg \text{UndirectedPath}(X, Y)$ must be attacked by Σ .

If A is $\text{UndirectedPath}(X, Y)$ and is not attacked, then Lemma 18 guarantees that $\neg \text{UndirectedPath}(X, Y)$ is also not in Σ and not attacked by it, which as we just saw, is impossible. Thus $\text{UndirectedPath}(X, Y)$ must be attacked by Σ .

Proving that $\text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ and $\neg \text{UndirectedPath}(X, Y) \text{ Excluding}(Z, W)$ must be attacked by Σ if they are not in Σ is proved in the same way as for $\text{UndirectedPath}(X, Y)$ and $\neg \text{UndirectedPath}(X, Y)$. The case of $\neg \text{DirectedPath}(X, Y)$ can be proved in a similar manner to $\neg \text{UndirectedPath}(X, Y)$, except that instead of only item 17 to establish the existence of the set N , we need to see that in any case at least one of items 14 to 16 implies the existence of another $\neg \text{DirectedPath}(Z, W)$ not in Σ and not attacked by Σ . That $\text{DirectedPath}(X, Y)$ is attacked then once again follows from Lemma 18. \square

Given this result we can thus efficiently identify a set Σ as being a preferred extensions of \mathcal{A}^g by simply checking if it attacks all arguments in $L \setminus \Sigma$ and only those.

6. Fusing agoras

We now address the problem of having agents agree on a preferred extension of \mathcal{A}^g , given that each of them has its own prior beliefs, as expressed by the compromise score function s_i , and that each know the combination function c . There has not been a lot of work done in dialectics for more than two agents, where the simple proponent/opponent dualism does not suffice. The solution that we propose here is inspired by the Risk Agoras of [15] and [16] and the traditional blackboard architecture of MAS of cooperating agents, without being an actual instantiation of any of them. We construct a *fusing agora*, which is a framework in which the agents can debate. The agora has the property that, if agents are allowed to run the debate to conclusion, they end up with the best possible compromise according to their joint compromise score, and that throughout the debate they maintain a compromise, which improves as the debate progresses.

In the agora we shall take a \mathcal{A}^g -candidate (I, O, U) as a unique representative of a partial compromise (I, O) . This is possible since I and O are subsets of L , and thus both contain sentences describing aspects of a compromise graph as required, and since U is determined by I and O . Any leaf candidate representing a preferred extension then uniquely identifies a completed pattern, as guaranteed by Theorem 16. Agents can explore all compromises by examining a $(\emptyset, \emptyset, L)$ -tree. Continually the agents take it upon themselves to explore sub-trees of this tree, and mark other sub-trees as open for investigation by other agents. The heuristics guiding the agents' choices for exploration, in addition to s_1, \dots, s_k and c , then determine the outcome.

The agora can work in a variety of ways, depending on the behavior of the individual agents (a vanilla algorithm for an individual agent is provided later in Algorithm 1), but basically builds on two elements, which we assume each agent can access in a synchronized fashion only: A *pool of candidates* C and a *current best result* $\langle I_*, s_{I_*} \rangle$. C consists of pairs $\langle C, s \rangle$, where C is an \mathcal{A}^g -candidate and thus a sub-tree of a $(\emptyset, \emptyset, L)$ -tree, and s is a real value. I_* is either the empty set or a preferred extension of \mathcal{A}^g , and s_{I_*} is a real value. Initially, C contains only one element $\langle (\emptyset, \emptyset, L), 0 \rangle$, and $\langle I_*, s_{I_*} \rangle$ is $\langle \emptyset, -\infty \rangle$.

Each agent i can utter the following locutions:

- *ExploreFromPool_i*($\langle C, s \rangle$)—where $\langle C, s \rangle$ is a member of C . The meaning of the locution is that agent i takes upon itself the responsibility to investigate the preferred extensions in a C -tree, assuming that C has a joint compromise score of s .
- *PutInPool_i*($\langle C, s \rangle$)—where C is an \mathcal{A}^g -candidate, and s is a real value. The meaning of the locution is that agent i wants someone else to investigate the preferred extensions in a C -tree, and that C has a joint compromise score of s .
- *UpdateBest_i*($\langle I, s \rangle$)—where I is a subset of L , and s is a real value. The meaning of the locution is that agent i has identified a preferred extension I with a joint compromise score s higher than s_{I^*} .
- *AskOpinion_i*(C_1, C_2)—where C_1 and C_2 are \mathcal{A}^g -candidates. The meaning of the locution is that agent i needs to know $s_j(C_1, C_2)$ for all other agents j .
- *StateOpinion_i*(C_1, C_2, s_δ)—where C_1 and C_2 are \mathcal{A}^g -candidates, and s_δ is a real value. The meaning of the locution is that $s_i(C_1, C_2)$ is s_δ .

The rules governing which locutions individual agents can utter, as well as their effects, we present as a set of pre and post conditions:

- *ExploreFromPool_i*($\langle C, s \rangle$)
 - Pre: $\langle C, s \rangle$ is in C .
 - Post: $\langle C, s \rangle$ is removed from C
- *PutInPool_i*($\langle C, s \rangle$)
 - Pre: There is no $\langle C', s' \rangle$ in C such that C is a sub-tree of some C' -tree.
 - Post: $\langle C, s \rangle$ is in C .
- *UpdateBest_i*($\langle I, s \rangle$)
 - Pre: $s > s_{I^*}$.
 - Post: $\langle I^*, s_{I^*} \rangle$ is set to $\langle I, s \rangle$.

Locutions *AskOpinion_i*() and *StateOpinion_i*() have no pre or post conditions attached.

The basic algorithm in Algorithm 1 corresponds to an exhaustive search, if it is followed by all agents. The search is gradual in two senses: One, the longer the search goes on, the more elements the average candidate in C will have in its I and O sets, and thus the closer it will be to describing a full compromise. Two, the current compromise held in I^* will have an increasingly higher score.

In order for the search to be a success, each agent i would of course need to keep an eye out for *AskOpinion_j*(\cdot)'s uttered by other agents, and reply to these with *StateOpinion_i*(\cdot). It is relatively easy to verify that agents using Algorithm 1 are uttering locutions in accordance with the pre and post conditions of the fusing agora. Algorithm 1 uses a series of helper functions, which are described below.

Prune($C \equiv (I, O, U)$) uses pruning rules to investigate whether there is an argument A in U such that either $C + A$ or $C - A$ contains no leaves with preferred extensions. If this is the case, the method invokes itself recursively on the sub-tree that did not get pruned away, until no further branches can be pruned. Some general pruning rules are given in [18], and more can be established for the specific case of \mathcal{A}^g . For instance, it is known that any preferred extension of \mathcal{A}^g is stable, so whenever $I \cup U$ does not attack an argument $A \in U$, then the branch corresponding to $C - A$ cannot contain any preferred extensions.

SelectCandidate(C) picks a promising candidate from C . A promising candidate could be one with a high score annotated, since these encode good partial compromises, or candidates with small U sets, as these represent partial compromises that are nearly complete. If all agents use the same criteria for picking promising candidates, this selection can be sped up by implementing the pool as a sorted list. *SelectCandidate*(\cdot) is one of the areas where heuristics limiting the search space can be implemented. For instance, it makes sense to allow agents to abstain from exploring the sub-tree rooted at a candidate if it cannot contain compromises that are consistent with their own BN. This would mean that in cases where agents agree on all or most of the aspects of \mathcal{G}_* only few candidates would need to be explored.

PreferredExtension(I) is a Boolean valued function that returns true if the conflict-free set I is a preferred extension of \mathcal{A}^g . The task of answering this is simplified by Theorem 19, as it states that I is a preferred extension iff I attacks each argument in $L \setminus I$.

```

1:  $\langle C, s \rangle \leftarrow \text{SelectCandidate}(C)$ 
2:  $\text{ExploreFromPool}_i(\langle C, s \rangle)$ 
3:  $C' \triangleq (I', O', U') \leftarrow \text{Prune}(C)$ 
4: if  $U' = \emptyset$  then
5:   if  $\text{PreferredExtension}(I')$  then
6:      $\text{AskOpinion}_i(C, C')$ 
7:      $s_i \leftarrow s_i(C, C')$ 
8:     wait for  $\text{StateOpinion}_j(C, C', s_j) \forall j \neq i$ 
9:      $s' \leftarrow c(s_1, \dots, s_k) + s$ 
10:    if  $s' > s_{I^*}$  then
11:       $\text{UpdateBest}_i(\langle C', s' \rangle)$ 
12:    end if
13:  end if
14:  go to 1
15: else
16:    $A \leftarrow \text{SelectArgument}(C')$ 
17:    $\text{AskOpinion}_i(C, C' + A)$ 
18:    $\text{AskOpinion}_i(C, C' - A)$ 
19:    $s_i^+ \leftarrow s_i(C, C' + A)$ 
20:    $s_i^- \leftarrow s_i(C, C' - A)$ 
21:   wait for  $\text{StateOpinion}_j(C, C' + A, s_j^+) \text{ and } \text{StateOpinion}_j(C, C' - A, s_j^-) \forall j \neq i$ 
22:    $s^+ \leftarrow c(s_1^+, \dots, s_k^+)$ 
23:    $s^- \leftarrow c(s_1^-, \dots, s_k^-)$ 
24:   if  $s^+ > s^-$  then
25:      $\text{PutInPool}_i(\langle C' - A, s + s^- \rangle)$ 
26:      $C \leftarrow C' + A$ 
27:      $s \leftarrow s + s^+$ 
28:   else
29:      $\text{PutInPool}_i(\langle C' + A, s + s^+ \rangle)$ 
30:      $C \leftarrow C' - A$ 
31:      $s \leftarrow s + s^-$ 
32:   end if
33:   go to 3
34: end if

```

Algorithm 1. Vanilla algorithm for agent i .

$\text{SelectArgument}(C \equiv (I, O, U))$ simply selects an element A of U . This selection can be based on the agent's own score increase going from C to $C + A$ or $C - A$, or it might involve negotiations or argument with other agents.

We illustrate the basics of the fusing agora and Algorithm 1 with an example.

Example 20. Refer back to Examples 2 and 3, and consider the simple case where the compromise scoring function of the three agents is defined as the number of features that correspond to the completed pattern of their model minus those that do not, as exemplified in Example 4. We shall follow the initial actions in the agora, where C consists of only one candidate, $\langle C \equiv (\emptyset, \emptyset, L), 0 \rangle$. We assume that Agent 1 is first to make a move.

Agent 1's first choice for exploration is simple, as C only contains one element, and the agent thus utters $\text{ExploreFromPool}_1(\langle C, 0 \rangle)$. This candidate cannot be pruned, and since the third set in it is not empty, Agent 1 selects an argument to add to one of the other two sets, say, $\text{Arc}(A, D)$. Agent 1 then constructs the two candidates (the third set of the candidates are left out, as it is simply L minus the first two sets):

$$\begin{aligned} \mathcal{C} + \text{Arc}(A, D) = & (\text{Arc}(A, D), \{\text{Arc}(D, A), \text{Link}(A, D), \\ & \text{NonAdjacent}(A, D), \text{ArcNotAllowed}(A, D), \\ & \neg \text{DirectedPath}(A, D), \text{DirectedPath}(D, A), \\ & \text{UndirectedPath}(A, D), \text{UndirectedPath}(D, A)\}) \end{aligned}$$

and

$$\mathcal{C} - \text{Arc}(A, D) = (\emptyset, \{\text{Arc}(A, D)\}).$$

Agent 1 ask the other two agents their opinion by $\text{AskOpinion}_1(\mathcal{C}, \mathcal{C} + \text{Arc}(A, E))$ and $\text{AskOpinion}_1(\mathcal{C}, \mathcal{C} - \text{Arc}(A, E))$, and receives four answers -1 , -1 , 1 , and 1 (as none of the other agents agree with $\text{Arc}(A, D)$). Added to its own perception, Agent 1 gets the scores $s^+ = -3$ and $s^- = 3$, and it thus utters $\text{PutInPool}_i((\mathcal{C} + \text{Arc}(A, E), -3))$, and continues exploring $\mathcal{C} - \text{Arc}(A, E)$.

Let us switch to Agent 2, who at this point utters $\text{ExploreFromPool}_2((\mathcal{C} + \text{Arc}(A, E) \equiv (\mathbf{I}_+, \mathbf{O}_+, \mathbf{U}_+), -3))$, and takes on the responsibility of examining the sub-tree rooted at $\mathcal{C} + \text{Arc}(A, E)$. First, Agent 2 prunes the candidate, by adding $\text{DirectedPath}(A, D)$ to \mathbf{I}_+ , and then chooses a new argument to investigate (seeing that \mathbf{U}_+ is not empty). For sake of coherence with previous examples, assume this argument is $\text{Arc}(D, E)$. Agent 2 creates the candidates

$$\begin{aligned} \mathcal{C} + \text{Arc}(A, D) + \text{Arc}(D, E) = & (\mathbf{I}_+ \cup \{\text{DirectedPath}(A, D), \text{Arc}(D, E)\}, \\ & \mathbf{O}_+ \cup \{\text{Arc}(E, D), \text{Link}(D, E), \text{NonAdjacent}(D, E), \\ & \text{ArcNotAllowed}(D, E), \neg \text{DirectedPath}(D, E), \\ & \text{DirectedPath}(E, D), \text{UndirectedPath}(D, E), \\ & \text{UndirectedPath}(E, D), \text{ArcNotAllowed}(A, E)\}) \end{aligned}$$

and

$$\mathcal{C} + \text{Arc}(A, D) - \text{Arc}(D, E) = (\mathbf{I}_+, \mathbf{O}_+ \cup \{\text{Arc}(D, E)\})$$

and gets an evaluation from Agents 1 and 3 using $\text{AskOpinion}_2(\cdot)$. These evaluations are $+1$, -1 , -1 , and $+1$, and taken together with $+1$ and -1 derived from Agent 2's own beliefs, it gets $s^+ = 1$ and $s^- = -1$, so it puts $\mathcal{C} + \text{Arc}(A, D) - \text{Arc}(D, E)$ in \mathcal{C} and continues exploring $\mathcal{C} + \text{Arc}(A, D) + \text{Arc}(D, E)$.

Of course, the debate in the agora can be stopped at any time, and $\mathcal{G}[\mathbf{I}_*^g]$ will then be the best compromise encountered so far, as it is only ever replaced by compromises having a higher joint compromise score. More specifically we have that:

Proposition 21. *Let agents 1 to k argue in a fusing agora. If*

- *all agents have used Algorithm 1, and*
- *each agent i have replied to all $\text{AskOpinion}_j(\cdot)$'s uttered by other agents, with a $\text{StateOpinion}_i(\cdot)$ consistent with s_i ,*

then all locutions uttered by agents are in accordance with the pre and post conditions of the fusing agora, and $(\mathbf{I}_, \mathbf{L} \setminus \mathbf{I}_*, \emptyset)$ has a higher joint compromise score s_i than all other explored leaf candidates $(\mathbf{I}, \mathbf{L} \setminus \mathbf{I}, \emptyset)$, where \mathbf{I} is a preferred extension of \mathcal{A}^g .*

If we furthermore have that

- *all agents have completed processing,*
- *the pool of candidates has not been thinned along the way, and*
- *the pool is empty now,*

then $\mathcal{G}[\mathbf{I}_^g] = \mathcal{G}_*$.*

It is worth stressing that Algorithm 1 is a vanilla algorithm, and that the agora is open for more aggressive behavior. One such behavior could be to have agents skip the asking for opinions part in lines 14 to 22 for most additions of arguments (and basing the decision only on the agents own beliefs), and only ask when the agent itself is indifferent. Another behavior could be to never perform lines 23 and 27, which would correspond to a myopic greedy construction of the compromise. Alternatively, these two lines could be carried out only when the difference between s^+ and s^- is very small. We could even have setups where the agents show different behaviors, or where individual agents change behavior during debate depending on their available resources and utility of a good compromise. Moreover, the agora does not require that agents wait for a candidate to be in the pool, before somebody can start exploring this candidate; so even when one agent is pursuing an aggressive strategy and fails to leave candidates for others to explore, other agents can still decide to explore these. The point is, that no matter what behavior is required, the basics of the agora and the agents remains the same, and can be reused.

7. Discussion

We have introduced a problem which we believe is a challenging one for the argumentation community, due to its mix of complexity and conditional decomposability as well as its origin in conflicting knowledge bases. Our own solution enables agents to judge the possible compromises resulting from a partial compromise, by constructing a candidate tree rooted in this partial compromise, and the agora we have proposed ensures that such exploration can take place in a distributed fashion. As an aside, we note that the method given here provides for a distributed solution to the problem of enumerating all equivalence classes of BNs [9].

One problem with the vanilla algorithm we have given, is that agents exploring a branch of a candidate-tree can end up putting a lot of candidates into the pool of annotated candidates. The space requirements for storing the pool of annotated candidates can be prohibitive, so it might be required that the candidates in the pool are defined in relation to each other, which imposes restrictions on which candidates an agent can choose to explore, as these are removed from the pool. Furthermore, it might be necessary to construct heuristics for thinning the pool of annotated candidates. These issues, as well as finding good heuristics for selecting candidates to explore are challenging topics for future research.

Acknowledgements

This work was partly supported by NSF REC-02-19347, NSF IIS-0329037, and EU PF6-IST 002307 (ASPIC).

References

- [1] S.A. Andersson, D. Madigan, M.D. Perlman, A characterization of Markov equivalence classes for acyclic digraphs, *Annals of Statistics* 25 (2) (1997) 505–541.
- [2] C. Cayrol, S. Doutre, J. Mengin, On decision problems related to the preferred semantics for argumentation frameworks, *Journal of Logic and Computation* 13 (3) (2003) 377–403.
- [3] D.M. Chickering, Learning equivalence classes of Bayesian-network structures, *Journal of Machine Learning Research* 2 (2002) 445–498.
- [4] J. Del Sagrado, S. Moral, Qualitative combination of Bayesian networks, *International Journal of Intelligent Systems* 18 (2) (2003) 237–249.
- [5] Y. Dimopoulos, B. Nebel, F. Toni, On the computational complexity of assumption-based argumentation for default reasoning, *Artificial Intelligence* 141 (1) (2002) 57–78.
- [6] S. Doutre, J. Mengin, Preferred extensions of argumentation frameworks: Query, answering and computation, in: R. Goré, A. Leitsch, T. Nipkow (Eds.), *Proceedings of the First International Joint Conference on Automated Reasoning*, in: *Lecture Notes in Computer Science*, vol. 2083, Springer-Verlag, 2001.
- [7] P.M. Dung, On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games, *Artificial Intelligence* 77 (2) (1995) 321–358.
- [8] M. Elvang-Goransson, A. Hunter, Argumentative logics: Reasoning with classically inconsistent information, *Data Knowledge Engineering* 16 (2) (1995) 125–145.
- [9] S. Gillispie, C. Lemieux, Enumerating Markov equivalence classes of acyclic digraph models, in: J. Breese, D. Koller (Eds.), *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann Publishers, 2001.
- [10] P.M.-R. II, U. Chajewska, Aggregating learned probabilistic beliefs, in: J. Breese, D. Koller (Eds.), *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann Publishers, 2001.
- [11] F.V. Jensen, *Bayesian Networks and Decision Graphs*, Springer-Verlag, 2001.
- [12] P. Krause, S. Ambler, M. Elvang-Goransson, J. Fox, A logic of argumentation for reasoning under uncertainty, *Computational Intelligence* 11 (1) (1995) 113–131.

- [13] S.L. Lauritzen, *Graphical Models*, Oxford University Press, 1996.
- [14] I. Matzkevich, B. Abramson, The topological fusion of Bayes nets, in: D. Dubois, M.P. Wellman, B. D'Ambrosio, P. Smets (Eds.), *Proceedings of the Eighth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann Publishers, 1992.
- [15] P. McBurney, S. Parsons, Risk agoras: Dialectical argumentation for scientific reasoning, in: C. Boutilier, M. Goldszmidt (Eds.), *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann Publishers, 2000.
- [16] P. McBurney, S. Parsons, Chance discovery using dialectical argumentation, in: T. Terano, T. Nishida, A. Namatame, S. Tsumoto, Y. Oh-sawa, T. Washio (Eds.), *New Frontiers in Artificial Intelligence*, in: Joint Japanese Society for Artificial Intelligence 2001 Workshop Post-Proceedings, vol. 2253, Springer-Verlag, 2001.
- [17] S.H. Nielsen, Reasoning and decision making with multiple autonomous agents, Ph.D. thesis, Aalborg University, 2007.
- [18] S.H. Nielsen, S. Parsons, Computing preferred extensions for argumentation systems with sets of attacking arguments, in: P.E. Dunne, T.J.M. Bench-Capon (Eds.), *Proceedings of the First International Conference on Computational Models of Argument*, in: *Frontiers in Artificial Intelligence and Applications*, vol. 144, IOS Press, 2006.
- [19] S.H. Nielsen, S. Parsons, A generalization of Dung's abstract framework for argumentation: AArguing with sets of attacking arguments, in: N. Maudet, S. Parsons, I. Rahwan (Eds.), *Proceedings of the Third Workshop on Argumentation in Multi-agent Systems*, Future University, Hakodate, Japan, 2006.
- [20] J. Pearl, *Probabilistic Reasoning in Intelligent Systems, Representation & Reasoning*, Morgan Kaufmann Publishers, 1988.
- [21] D.M. Pennock, M.P. Wellman, Graphical representations of consensus belief, in: K. Laskey, H. Prade (Eds.), *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann Publishers, 1999.
- [22] J.L. Pollock, *Cognitive Carpentry: A Blueprint for How to Build a Person*, MIT Press, 1995.
- [23] M. Richardson, P. Domingos, Learning with knowledge from multiple experts, in: T. Fawcett, N. Mishra (Eds.), *ICML 20*, AAAI Press, 2003.
- [24] S. Saha, S. Sen, A Bayes net approach to argumentation based negotiation, in: *ArgMAS 1*, Springer-Verlag, 2004.
- [25] G.R. Simari, R.P. Loui, A mathematical treatment of defeasible reasoning and its implementation, *Artificial Intelligence* 53 (1992) 125–157.
- [26] T. Verma, J. Pearl, Equivalence and synthesis of causal models, in: L.K.J.L.P. Bonissone, M. Henrion (Eds.), *Proceedings of the Sixth Conference on Uncertainty in Artificial Intelligence*, Elsevier Science Publishing, 1991.
- [27] G. Vreeswijk, Bayesian inference and defeasible reasoning: Suggestions for a unified framework, working paper, Department of Computer Science, University of Utrecht, 1999.
- [28] G. Vreeswijk, H. Prakken, Credulous and skeptical argument games for preferred semantics, in: M. Ojeda-Aciego, I.P. de Guzmán, G. Brewka, L.M. Pereira (Eds.), *Proceedings of the Seventh European Workshop on Logic in Artificial Intelligence*, in: *Lecture Notes in Computer Science*, vol. 1919, Springer-Verlag, 2000.
- [29] G.A.W. Vreeswijk, Abstract argumentation systems, *Artificial Intelligence* 90 (1) (1997) 225–279.
- [30] G.A.W. Vreeswijk, Argumentation in Bayesian belief networks, in: *ArgMAS 1*, in: *Lecture Notes in Artificial Intelligence*, Springer-Verlag, 2004.