

The Carneades Argumentation Framework – Using Presumptions and Exceptions to Model Critical Questions

Thomas F. Gordon¹ and Douglas Walton²

Abstract. In 2005, Gordon and Walton presented initial ideas for a computational model of defeasible argument [12, 26], which builds on and elaborates Walton’s theory of argumentation [28, 31]. The current paper reports on progress which has been made in the meantime. It presents a formal, mathematical model of argument evaluation which applies proof standards [8] to determine the defensibility of arguments and the acceptability of statements on an issue-by-issue basis. The main original contribution of the Carneades Argumentation Framework is its use of three kinds of premises (ordinary premises, presumptions and exceptions) and information about the dialectical status of statements (undisputed, at issue, accepted or rejected) to model critical questions in such a way as to allow the burden of proof to be allocated to the proponent or the respondent, as appropriate. Both of these elements are required for this purpose: presumptions hold without supporting argument only so long as they have not been put at issue by actually asking the critical question.

1 Introduction

The work in this paper flows from previous attempts to solve a key problem common to AI and argumentation theory concerning the using of the device of critical questions to evaluate an argument. Critical questions were first introduced by Arthur Hastings [15] as part of his analysis of presumptive argumentation schemes. The critical questions attached to an argumentation scheme enumerate ways of challenging arguments created using the scheme. The current method of evaluating an argument that fits a scheme, like that for argument from expert opinion, is by a shifting of the burden of proof from one side to the other in a dialog [30]. When the respondent asks one of the critical questions matching the scheme, the burden of proof shifts back to the proponent’s side, defeating or undercutting the argument until the critical question has been answered successfully. At least this has been the general approach of argumentation theory. Recently, however, it was observed [3] that critical questions differ with respect to their impact on the burden of proof. These observations led to two theories about the shifting of the burden of proof when critical questions are asked. According to one theory,

when any critical question is asked, the burden shifts to the proponent’s side to answer the question and, if no answer is given, the argument fails. According to the other theory, merely asking a critical question is not enough to shift the burden of proof back to the proponent. On this theory, to make the argument fail, the question needs to be supported by further argument. Some critical questions fit one theory better, while others fit the other theory better. This duality has posed a recurring problem for the project of formalizing schemes.

In this paper, we put forward a new model for evaluating defeasible arguments that solves this problem, continuing work we began in 2005 [12, 26]. The current paper presents a formal, mathematical model of argument evaluation which applies proof standards [8] to determine the defensibility of arguments and the acceptability of statements on an issue-by-issue basis. The formal model is called the Carneades Argumentation Framework, in honor of the Greek skeptic philosopher who emphasized the importance of plausible reasoning [6, vol. 1, p. 33-34].

Arguments in Carneades are identified, analyzed and evaluated not only by fitting premise-conclusion structures that can be identified using argumentation schemes. Arguments also have a dialectical aspect, in that they can be seen as having been put forward on one side or the other of an issue during a dialog. The evaluation of arguments in Carneades depends on the stage of the dialog. Whether or not a premise of an argument holds depends on whether it is undisputed, at issue, or decided. One way to raise an issue is to ask a critical question. Also, the proof standard applicable for some issue may depend on the stage of the dialog. In a deliberation dialog, for example, a weak burden of proof would seem appropriate during brainstorming, in an early phase of the dialog. The Carneades Argumentation Framework is designed to be used in a layered model of dialectical argument [19] for various kinds of dialogs, where higher layers are responsible for modeling such things as speech acts, argumentation protocols and argument strategies.

The rest of the paper is structured as follows. The next two sections formally define the Carneades Argumentation Framework. Section 2 defines the structure of arguments and illustrates this structure with examples from related work by Toulmin, Pollock and others. Section 3 formally defines how arguments are evaluated in terms of the acceptability of statements, the defensibility of arguments, and the satisfiability

¹ Fraunhofer FOKUS, Berlin, Germany, email: thomas.gordon@fokus.fraunhofer.de

² Department of Philosophy, University of Winnipeg, Winnipeg, Manitoba, Canada, email: d.walton@uwinnipeg.ca

of proof standards. Section 4 illustrates argument evaluation with an example from the AI and Law literature. The paper closes in Section 5 with a brief discussion of related work and some ideas for future work.

2 Argument Structure

We begin by defining the structure of arguments. Unlike Dung’s model [5], in which the internal structure of arguments is irrelevant for the purpose of determining their defensibility, our model makes use of and depends on the more conventional conception of argument in the argumentation theory literature, in which arguments are a kind of conditional linking a set of premises to a conclusion. Intuitively, the premises and the conclusion of arguments are statements about the world, which may be accepted as being true or false. In [12] the internal structure of statements was defined in such a way as to enable the domain of discourse to be modeled in a way compatible with emerging standards of the Semantic Web [2]. These details, however, need not concern us here. For the purpose of evaluating arguments, the internal structure of statements is not important. We only require the ability to compare two statements to determine whether or not they are equal.

Definition 1 (Statements) Let $\langle \text{statement}, = \rangle$ be a structure, where *statement* denotes the set of declarative sentences in some language and $=$ is an equality relation, modeled as a function of type $\text{statement} \times \text{statement} \rightarrow \text{boolean}$.

Next, to support defeasible argumentation and allow the burden of proof to be distributed, we distinguish three kinds of premises.

Definition 2 (Premises) Let *premise* denote the set of premises. There are three kinds of premises:

1. If s is a statement, then $\text{premise}(s)$ is a premise. These are called ordinary premises. As a notational convenience, we will use a statement s alone to denote $\text{premise}(s)$, when the context makes it clear that the statement is being used as a premise.
2. If s is a statement, then $\bullet s$, called a presumption, is a premise.
3. If s is a statement, then $\circ s$, called an exception, is a premise.
4. Nothing else is a premise.

Now we are ready to define the structure of arguments.

Definition 3 (Arguments) An argument is a tuple $\langle c, d, p \rangle$, where c is a statement, $d \in \{\text{pro}, \text{con}\}$ and $p \in \mathcal{P}(\text{premise})$. If a is an argument $\langle c, d, p \rangle$, then $\text{conclusion}(a) = c$, $\text{direction}(a) = d$ and $\text{premises}(a) = p$. Where convenient, *pro* arguments will be notated as $p_1, \dots, p_n \rightarrow c$ and *con* arguments as $p_1, \dots, p_n \dashv\!\!\!\dashv c$.

This approach, with two kinds of arguments, *pro* and *con*, is somewhat different than the argument diagramming model developed by Walton in [28] and implemented in Araucaria. There counterarguments are modelled as arguments *pro* some statement which has been asserted to be in conflict with the

conclusion of the other argument, called a *refutation*. Our approach, with its two kinds of arguments, is not uncommon in the literature on defeasible argument [18, 22, 14, 13].

We assume arguments are asserted by the participants of a dialog. We have specified and implemented a simple communication language and argumentation protocol to test Carneades, but that is a subject for another paper. For our purposes here, it is sufficient to note that argument moves, i.e. speech acts, are modelled as functions which map a state of the dialog to another state. (Again, this is a purely functional model, so states are not modified.) A dialog state is a tuple $\langle t, h, G \rangle$, where t is a statement, the *thesis* of the dialog, h is a sequence of moves, representing the history of the dialog, and G is an *argument graph*.³

It is these argument graphs which concern us here. An argument graph plays a role comparable to a set of formulas in logic. Whereas in logic the truth of a formula is defined in terms of a (consequence) relation between sets of formulas, here we will define the *acceptability* of statements in argument graphs. An argument graph is not merely a set of arguments. Rather, as its name suggests, it is a finite graph. There are two kinds of nodes, statement nodes and argument nodes. The edges of the graph link up the premises and conclusions of the arguments. Each statement is represented by at most one node in the graph.

To illustrate argument graphs, suppose we have the following (construed) arguments from the domain of contract law:

- a1. agreement, \circ minor \rightarrow contract
- a2. oral, \bullet estate $\dashv\!\!\!\dashv$ contract
- a3. email \rightarrow oral
- a4. deed $\dashv\!\!\!\dashv$ agreement
- a5. \bullet deed \rightarrow estate

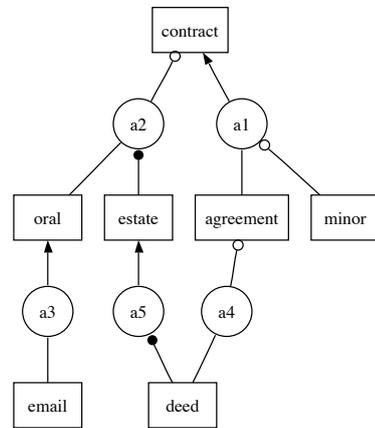


Figure 1. Argument Graph

The argument graph induced by these arguments is shown in Figure 1. In this figure, statements are displayed as boxes and arguments as circles. Different arrowhead shapes are used to distinguish *pro* and *con* arguments as well as the three

³ In prior work [11, 13], Gordon has referred to argument graphs as *dialectical graphs*.

kinds of premises. Pro arguments are indicated using ordinary arrowheads; con arguments with open-dot arrowheads. Ordinary premises are represented as edges with no arrowheads, presumptions with closed-dot arrowheads and exceptions with open-dot arrowheads. (The direction of the edge is implicit in the case of ordinary premises; the direction is always from the premise to the argument.) Notice that the premise type cannot be adequately represented using statement labels, since argument graphs are not restricted to trees. A statement may be used in multiple arguments and as a different type of premise in each argument. The above example illustrates this point. The fourth and the fifth arguments each use the statement ‘deed’ in a premise. In the fourth argument it is used in an ordinary premise but in the fifth it is used in a presumption. Walton has called this use of shared premises a *divergent argument structure* [28, p. 91].

Although argument graphs are not restricted to trees, they are not completely general; we do not allow cycles. This restriction assures the decidability of the defensibility and acceptability properties of arguments and statements, respectively.

Definition 4 (Argument Graphs) *An argument-graph is a labeled, finite, directed, acyclic, bipartite graph, consisting of argument nodes and statement nodes. The edges link the argument nodes to the statements in the premises and conclusion of each argument.*

This completes the formal definition of the structure of arguments and argument graphs. Let us now discuss briefly the expressiveness of this model, beginning by comparing our approach with Toulmin’s model [21]. Recall that arguments in Toulmin’s model consist of a single premise, called the *datum*; a conclusion, called the *claim*; a kind of rule, called the *warrant*, which supports the inference from the premise to the conclusion of the argument; an additional piece of data, called *backing*, which provides support for the warrant; an exception, called a *rebuttal*; and, finally, a *qualifier* stating the probative value of the inference (e.g. presumably, or necessarily). Of these, the datum and conclusion are handled in a straightforward way in our model. The set of premises of an argument generalizes the single datum in Toulmin’s system. Claims are modeled comparably, as conclusions. Rebuttals are modeled with con arguments. The probative weight of an argument is handled as part of our model of proof standards, as will be explained shortly.

This leaves our interpretation of warrants and backing to be explained. Our model does not directly allow arguments about other arguments. (The conclusion of an argument must be a statement.) Rather, the approach we prefer is to add a presumption for the warrant to the premises of an argument. If an argument does not have such a presumption, the argument graph can first be extended to add one. We leave it up to the argumentation protocol of the procedural model to regulate under what conditions such *hidden premises* may be *revealed*. In effect, the datum and warrant are modelled as minor and major premises, much as in the classical theory of syllogism. Backing, in turn, can be modelled as a premise of an argument supporting the warrant.

For example, here is a version of Toulmin’s standard example about British citizenship.

Datum. Harry was born in Bermuda.

Claim. Harry is a British subject.

Warrant. A man born in Bermuda will generally be a British subject.

Backing. Civil Code §123 provides that persons born in Bermuda are generally British subjects.

Exception. Harry has become an American citizen.

The argument can be reconstructed in our framework as illustrated in Figure 2.

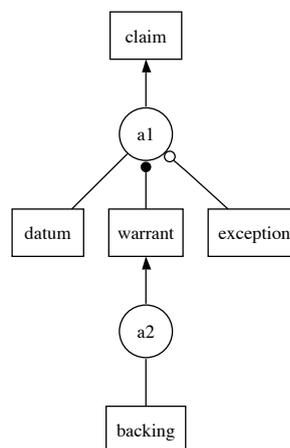


Figure 2. Reconstruction of Toulmin Diagrams

This approach generalizes Toulmin’s model, by supporting arguments pro and contra both warrants and backing, using the same argumentation framework as for arguments about any other kind of claim. Indeed, Toulmin appears to have overlooked the possibility of arguing against warrants or making an issue out of backing claims.

Our model of argument is rich enough to handle Pollock’s concepts of rebuttal, premise defeat and undercutting defeaters [18]. Rebuttals can be modeled as arguments in the opposite direction for the same conclusion. (If an argument a_1 is *pro* some statement s , then some argument a_2 *con* s is a rebuttal of a_1 , and vice versa.) Premise defeat can be modeled with arguments con an ordinary premise or presumption, or pro an exception.

Undercutting defeaters are a bit trickier. The idea of an undercutting defeater is to argue against the argument itself, or the rule or warrant which was applied to create the argument. We model undercutting defeaters by revealing and then attacking premises, similar to the way we handled warrants in the reconstruction of Toulmin’s system. Consider Pollock’s example of things which look red but turn out to be illuminated by a red light:

Red. The object is red.

Looks Red. The object looks red.

Applicable. The general rule “Things which look red are red.” applies to this object.

Illuminated. The object is illuminated by a red light.

An argument graph for this example is shown in Figure 3. Rather than undercutting argument a_1 (the object is red

because it looks red) directly, with an argument contra a_1 , we undercut the argument by first revealing a presumption (about the general rule being applicable in this case) and then assert an argument contra this presumption. Notice by the way that another presumption is still implicit in this example, namely a presumption for the “warrant” about things which look red being red.

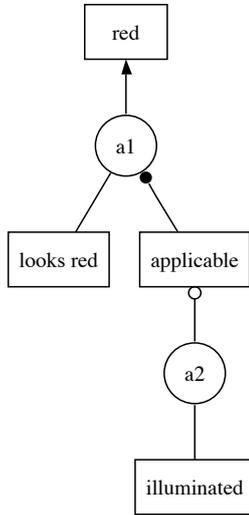


Figure 3. Undercutting Defeater Example

Walton [28] distinguishes two kinds of arguments, called *convergent* and *linked* arguments. Convergent arguments provide multiple reasons for a conclusion, each of which alone can be sufficient to accept the conclusion. Convergent arguments are handled in our approach by multiple arguments for the same conclusion. Linked arguments, on the other hand, consist of two or more premises which all must hold for the argument to provide significant support for its conclusion. Linked arguments are handled in our approach by defining arguments to consist of a set of premises, rather than a single premise, and defining arguments to be defensible only if all of their premises hold. (The concept of argument defensibility is formally defined below.)

Presumptions and exceptions are a refinement of Walton’s concept of *critical questions* [29]. Critical questions enumerate specific ways to defeat arguments matching some argument scheme. But so long as an issue has not been raised by actually asking some critical question, we would like to be able to express which answer to presume. The distinction between presumptions and exceptions here provides this ability.

Consider the scheme for arguments from expert opinion [25]:

Major Premise. Source E is an expert in the subject domain S containing proposition A .

Minor Premise. E asserts that proposition A in domain S is true.

Conclusion. A may plausibly be taken as true.

The scheme includes six critical questions:

- CQ1.** How credible is E as an expert source?
- CQ2.** Is E an expert in the field that A is in?
- CQ3.** Does E ’s testimony imply A ?
- CQ4.** Is E reliable?
- CQ5.** Is A consistent with the testimony of other experts?
- CQ6.** Is A supported by evidence?

When the scheme for arguments from expert opinion is instantiated to create a specific argument, the critical questions can be represented, in our model, as presumptions and exceptions. Whether a presumption or exception is appropriate depends on the burden of proof. If the respondent, the person who poses the critical question, should have the burden of proof, then the critical question should be modeled as an exception. If, on the other hand, the proponent, the party who used the schema to construct the argument, should have the burden of proof, then the critical question should be modeled as a presumption.⁴

Our model does not require that premises for critical questions be made explicit at the time the argument is first made. Rather, they can be *revealed* incrementally during the course of the dialog. The conditions under which a premise may be left implicit or revealed raise procedural issues which need to be addressed in the protocol for the type of dialog. Our contribution here is to provide an argumentation framework which can be used for modeling such protocols.

3 Argument Evaluation

By argument evaluation we mean determining whether a statement is *acceptable* in an argument graph. As we will see soon, this in turn will depend on the *defensibility* of arguments in the graph. Notice that our terminology is somewhat different than Dung’s [5], who speaks of the acceptability of arguments, rather than their defensibility. Also, for those readers familiar with our preliminary work on this subject in [12], please notice that the terminology and other details of the current model are different, even though the basic ideas and general approach are quite similar.

The definition of the acceptability of statements is recursive. The acceptability of a statement depends on its *proof standard*. Whether or not a statement’s proof standard is *satisfied* depends on the defensibility of the arguments pro and con this statement. The defensibility of an argument depends on whether or not its premises *hold*. Finally, we end up where we began: Whether or not a premise holds can depend on whether or not the premise’s statement is acceptable. Since the definitions are recursive, we cannot avoid making forward references to functions which will be defined later.

To evaluate a set of arguments in an argument graph, we require some additional information. Firstly, we need to know the current *status* of each statement in the dialog, i.e. whether it is accepted, rejected, at issue or undisputed. This status information is pragmatic; the status of statements is set by speech acts in the dialog, such as asking a question, asserting an argument or making a decision. Secondly, we assume that a proof standard has been assigned to each statement. We do

⁴ We agree with Verheij [24] that critical questions which are entailed by the premises of the argument schema are redundant and may be omitted. This is arguably the case in the example for the first three critical questions.

not address the question of how this is done. Presumably this will depend on domain knowledge and the type of dialog. Finally, one of the proof standards we will define, *preponderance of the evidence*, makes use of numerical weights, comparable to conditional probabilities. To use this proof standard, we require a weighing function.

Let us formalize these requirements by postulating an *argument context* as follows.

Definition 5 (Argument Context) Let C , the argument context, be a tuple $\langle G, \text{status}, \text{proof-standard}, \text{weight} \rangle$, where G is an argument-graph, status is a function of type $\text{statement} \rightarrow \{\text{accepted}, \text{rejected}, \text{undisputed}, \text{issue}\}$, proof-standard is a function of type $\text{statement} \rightarrow \{\text{SE}, \text{PE}, \text{DV}, \text{BRD}\}$ and weight is a function of type $\text{statement} \times \text{statement} \rightarrow \{0, \dots, 10\}$

Intuitively, a statement which has been used in a dialog is initially undisputed. Later in the dialog, an issue can be made out of this statement. Presumably after arguments pro and con have been collected for some period of time, a decision will be taken and the statement will be either accepted or rejected. The details of how this is done need not concern us further here. These are matters which need to be addressed fully when modeling protocols for dialogs.

Definition 6 (Acceptability of Statements)

Let acceptable be a function of type $\text{statement} \times \text{argument-graph} \rightarrow \text{boolean}$. A statement is acceptable in an argument graph if and only if it satisfies its proof standard in the argument graph: $\text{acceptable}(s, ag) = \text{satisfies}(s, \text{proof-standard}(s), ag)$.

Definition 7 (Satisfaction of Proof Standards)

A proof standard is a function of type $\text{statement} \times \text{argument-graph} \rightarrow \text{boolean}$. Let f be a proof standard. $\text{satisfies}(s, f, G) = f(s, G)$

Four proof standards are defined in this paper.

SE. A statement meets this standard iff it is supported by at least one defensible pro argument.

PE. A statement meets this standard iff its strongest defensible pro argument outweighs its strongest defensible con argument. This standard balances arguments using probabilistic weights.

DV. A statement meets this standard iff it is supported by at least one defensible pro argument and none of its con arguments are defensible.

BRD. A statement meets this standard iff it is supported by at least one defensible pro argument, all of its pro arguments are defensible and none of its con arguments are defensible.

The names of three of these standards are meant to suggest three legal proof standards: scintilla of evidence, preponderance of the evidence and beyond a reasonable doubt. However, we do not claim that the definitions of these standards, above, fully capture their legal meanings. What these standards have in common with their legal counterparts is their relative strength. If a statement satisfies a proof standard, it will also satisfy all weaker proof standards.

The name of the DV proof standard is an acronym for *dialectical validity*, a term used by Freeman and Farley [8]. They defined five proof standards. In addition to the four we have defined here, they included a fifth, called *beyond a doubt*, which was defined to be an even stronger standard than *beyond a reasonable doubt*.

The preponderance of evidence (PE) standard compares the weight of arguments. The weight of an argument is defined to be the same as the weight of its *weakest premise*, i.e., to be precise, the same as the weight of the premise with the lowest weight. Recall we assume a weighing function, weight , as part of the context to provide this information. The weight of a premise p for a conclusion c is $\text{weight}(p, c)$. Other proof standards which aggregate and compare weights are conceivable. For example, one could sum the weights of the arguments pro and con and compare these sums.

We have defined weights to be natural numbers in the range of 0 to 10. We originally considered using real numbers in the range of 0.0 to 1.0, as in probability theory. However, on the assumption that the weights will be estimated by human users, we prefer to use a simpler ordinal scale, since we are skeptical that users can estimate such weights with a greater degree of accuracy.

All of the proof standards defined above depend on a determination of the *defensibility* of arguments. Defensibility is defined next.

Definition 8 (Defensibility of Arguments)

Let defensible be a function of type $\text{argument} \times \text{argument-graph} \rightarrow \text{boolean}$. An argument α is defensible in an argument graph G if and only if all of its premises hold in the argument graph: $\text{defensible}(\alpha, G) = \text{all}(\lambda p. \text{holds}(p, G))(\text{premises } \alpha)$.⁵

Finally, we come to the last definition required for evaluating arguments, for the holds predicate. This is where the status of a statement in the argument context and the distinction between ordinary premises, presumptions and exceptions come into play. Accepted presumptions and ordinary premises hold. Rejected presumptions and ordinary premises do not hold. Undisputed presumptions hold. Undisputed ordinary premises do not hold. An exception, $\bullet s$, holds only if premise(s) does not hold.

Definition 9 (Holding of Premises) Let holds be a function of type $\text{premise} \times \text{argument-graph} \rightarrow \text{boolean}$. Let $\sigma = \text{status}(s)$. Whether or not a premise holds depends on its type (ordinary, presumption, or exception). Thus, there are the following three cases:

If p is an ordinary premise, $\text{premise}(s)$, then

$$\text{holds}(p, G) = \begin{cases} \text{true} & \text{if } \sigma = \text{accepted} \\ \text{false} & \text{if } \sigma = \text{rejected} \\ \text{acceptable}(s, G) & \text{if } \sigma = \text{issue} \\ \text{false} & \text{if } \sigma = \text{undisputed} \end{cases}$$

If p is a presumption, $\bullet s$, then

⁵ Here ‘all’ is a higher-order function, not a quantifier, applied to an anonymous function, represented with λ , as in lambda calculus.

$$\text{holds}(p, G) = \begin{cases} \text{true} & \text{if } \sigma = \text{accepted} \\ \text{false} & \text{if } \sigma = \text{rejected} \\ \text{acceptable}(s, G) & \text{if } \sigma = \text{issue} \\ \text{true} & \text{if } \sigma = \text{undisputed} \end{cases}$$

Finally, if p is an exception, $\circ s$, then

$$\text{holds}(p, G) = \neg \text{holds}(\text{premise}(s), G)$$

The important thing to notice is that whether or not a premise holds depends in this model not only on the arguments which have been asserted, but also on the kind of premise (ordinary, presumption, or exception) and the status of the premise’s statement in the argument graph (undisputed, at issue, accepted, or rejected). We assume that the status of a statement progresses in the course of the dialog:

1. Initially, statements used in arguments are undisputed. Whether or not a premise which uses this statement holds at this stage of the dialog depends on the kind of premise. Ordinary premises do not hold; presumptions do hold. *This is the only semantic difference between ordinary premises and presumptions in our model.* An exception holds at this stage only if it would not hold if it were an *ordinary premise*. Notice that exceptions are not the dual of presumptions. As undisputed presumptions hold, an undisputed exception would not hold if we had defined exceptions to hold only if they would not hold if they were presumptions. But this is not the semantics we want. Rather, both undisputed exceptions and undisputed presumptions hold.
2. At some point a participant may make an issue out of a statement. Now ordinary premises and presumptions which use this statement hold only if they are acceptable, i.e. only if the statement meets its proof standard, given the arguments which have been asserted. Exceptions at issue hold only if the statement is not acceptable. We presume that arguments will be exchanged in a dialog for some period of time, and that during this phase the acceptability of statements at issue will be in flux.
3. Finally, at some point a decision will be made to either accept or reject some statement at issue. The model does not constrain the discretion of users to decide as they please. Unacceptable statements may be accepted and acceptable statements may be rejected. This remains transparent however. Any interested person can check whether the decisions are justified given the arguments made and the applicable proof standards. Anyway, after a decision has been made, it is respected by the model: Accepted statements hold and rejected statements do not hold, no matter what arguments have been made or what proof standards apply.

4 An Example

Although our model of argument is rather simple, we claim, it is nonetheless rather difficult to illustrate all of its features, or indeed validate the model, with just a few examples. We have rather ambitious aims for the model. It should be sufficient for use as the *argumentation framework* layer [19] in procedural models of protocols for a wide variety of dialog types [31]. It should be sufficient as a basis for formal models of argument

schemes, including critical questions. The distinction between the three kinds of premises should be adequate for allocating the burden of proof. It should be capable of being extended to handle other proof standards, such as more adequate models of legal proof standards. And of course it should yield intuitive results when applied to real examples of natural arguments. We have begun the work of testing and validating the model, but much work remains. Here we can only present a couple of examples to illustrate its main features.

As we are particularly interested in legal applications, we have reconstructed several examples from the Artificial Intelligence and Law literature [11, 17, 24, 1]. Some of these [11, 17] are procedural models of argumentation. Our reconstruction of these examples makes use of a procedural model of persuasion dialogs, based on the argumentation framework presented here. For lack of space, we will instead illustrate the model with one of the other examples which does do require us to address these procedural aspects.

We have selected one of Verheij’s main examples [24, p. 69], which he calls the “grievous bodily harm” example. The example consists of the following statements.

- 8 years.** The accused is punishable by up to 8 years in imprisonment.
- bodily harm rule.** Inflicting grievous bodily harm is punishable by up to 8 years imprisonment.
- Article 302.** According to article 302 of the Dutch criminal code, inflicting grievous bodily harm is punishable by up to 8 years imprisonment.
- bodily harm.** The accused has inflicted grievous bodily harm upon the victim.
- 10 witnesses.** 10 pub customers’ testimonies: the accused was involved in the fight.
- accused’s testimony** I was not involved in the fight.
- broken ribs not sufficient.** Several broken ribs do not amount to grievous bodily harm.
- precedent 1.** The rule that several broken ribs does not amount to grievous bodily harm, explains precedent 1.
- lex specialis.** The rule explaining precedent 2 is more specific than the rule explaining precedent 1.
- sufficient with complications.** Several broken ribs with complications amount to grievous bodily harm.
- precedent 2.** The rule that several broken ribs with complications amount to grievous bodily harm, explains precedent 2.
- hospital report.** The victim has several broken ribs, with complications.

The arguments are displayed, together with their evaluation, in Figure 4. We’ve made some assumptions about the context, for the purposes of illustration:

- The status of statements is indicated in the diagram via a suffix: A question mark (?) means the statement is at issue; A plus sign (+) means it has been accepted; a minus sign (-) indicates it has been rejected; and the lack of a suffix means the statement is undisputed. The *lex specialis* and *10 witnesses* statements have been accepted. The statements of other leaf nodes are undisputed. All the other statements are at issue.
- The DV proof standard (dialectical validity) applies to all statements. This is closest to the evaluation criteria of Ver-

heij’s model of argumentation, which does not support multiple proof standards.

- Weights are irrelevant in this example, since the PE proof standard (preponderance of the evidence) is not used.

Some further assumptions about the types of the premises have been made, to illustrate many features of the system with this one example. The result of the evaluation has been indicated in the diagram by filling in the nodes for acceptable statements and defensible arguments with a gray background. All the other statements are not acceptable and all other arguments are not defensible. Let us now try to explain the result, for each issue:

- The main issue, or thesis, that the accused is punishable by up to 8 years in prison, is acceptable. This is because both premises of the argument a_1 are acceptable and there are no rebuttals to consider.
- The statement about the *bodily harm rule* is acceptable, because it is supported by one defensible argument, a_2 , and there are no counterarguments. Argument a_2 is defensible, because its single premise, about Article 302, is an undisputed presumption.
- The claim that the accused has inflicted bodily harm is acceptable, because it is supported by a defensible argument, a_3 , and neither of the two counterarguments are defensible. The supporting argument, a_3 , is defensible because its premise has been accepted.
- Argument a_4 is not defensible, because its premise, regarding the accused’s testimony, in which he claims not to have been involved in a fight, is at issue and not acceptable.
- The accused’s testimony is not acceptable for two reasons: 1) it is successfully countered by the argument a_6 , with the testimony of 10 witnesses who claim to have seen the fight. (This testimony has been accepted with no further argument or evidence.) 2) It is not supported by at least one defensible pro argument, as required by the DV proof standard.
- The statement about broken ribs not being sufficient to amount to grievous bodily harm is not acceptable both because its only pro argument, a_9 , is not defensible and also because its counterargument, a_7 , is defensible. That is, the statement would not have met the DV proof standard even if its supporting argument had been defensible, since it is countered by a_7 .
- The statement about several broken ribs with complications being grievous bodily harm is acceptable, because it is supported by a defensible argument, a_8 , and has no counterarguments. The argument a_8 is defensible, because its only premise, about the second precedent, has been presumed and is not at issue.
- Finally, argument a_9 is not defensible, although it is supported by an undisputed premise, about the first precedent, because the *lex specialis* exception has been revealed (we assume) and accepted. Notice how *lex specialis*, which provides a reason to prefer precedent 2 over precedent 1, can be modeled even though our argumentation framework does not explicitly provide a way to order arguments.

One important function of an argumentation framework is to provide a basis for clear and comprehensible explanations or justifications of decisions. Argumentation framework which

depend on a deep understanding of mathematics (e.g. fixed points) or formal logic (e.g. entailment from minimal subsets of hypotheses, as in some models of abduction) for justifying decision do not meet this requirement. We hope the Carneades system is sufficiently simple that explanations, such as the above, can be quickly appreciated and understood by people with no formal background in logic or mathematics.

5 Discussion

The idea of developing a computer model for managing support and justification relationships between propositions goes back to research on “truth” or reason maintenance systems in Artificial Intelligence [4, 16]. The first author’s prior work on the Pleadings Game [11] included a formal model of dialectical graphs, for recording various kinds of support and defeat relationships among arguments. The concept of an *argumentation framework* was introduced by Henry Prakken [19] as part of a three-layered model for dialectical systems. As noted previously, Freeman and Farley [8] were the first to our knowledge to develop a computational model of burden of proof.

The Zeno Argumentation Framework [13] was based on Horst Rittel’s Issue-Based Information System (IBIS) model of argumentation [20]. The Carneades Argumentation Framework, in contrast, uses mainstream argumentation theory as its starting point. Also, Zeno did not provide a foundation for modeling argument schemes with critical questions, and was not as well suited as the current system for modeling persuasion dialogs.

Verheij’s work in [23] was the source of inspiration for distinguishing between different kinds of critical questions, which we have called presumptions and exceptions. Verheij’s book, *Virtual Arguments* [24], includes an enlightening comparison of several theories of defeasible argumentation. Verheij compared them with regard to whether and, if so, how each system modeled 1) pro and con arguments; 2) warrants, in Toulmin’s sense; 3) argument evaluation; and, finally 4) theory construction. We have already explained how our formal model handles the first three of these dimensions. In our model, the set of statements found to be acceptable can be viewed as a theory constructed collaboratively by participants in a dialog. Indeed, the first author, influenced by Fiedler [7], has long viewed reasoning explicitly as a theory construction process [9, 10] and was first attracted to argumentation theory precisely for this reason.

One key element of our theory construction approach is the idea of revealing hidden or implicit premises during a dialog. This approach was illustrated during the discussion of Toulmin and Pollock, for example, where warrants and undercutting defeaters were modelled as implicit presumptions revealed during dialog. Walton and Reed have done some recent work showing how argument schemes can be used to reveal implicit premises [27].

The formal model has been fully implemented, in a declarative way using a functional programming language, and tested on a number of examples from the Artificial Intelligence and Law literature, thus far yielding intuitively acceptable results. This validation work is continuing. More work is required to validate the models of the various proof standards, in particular the model of preponderance of the evidence, which uses weights. For this purpose, we plan to reconstruct examples

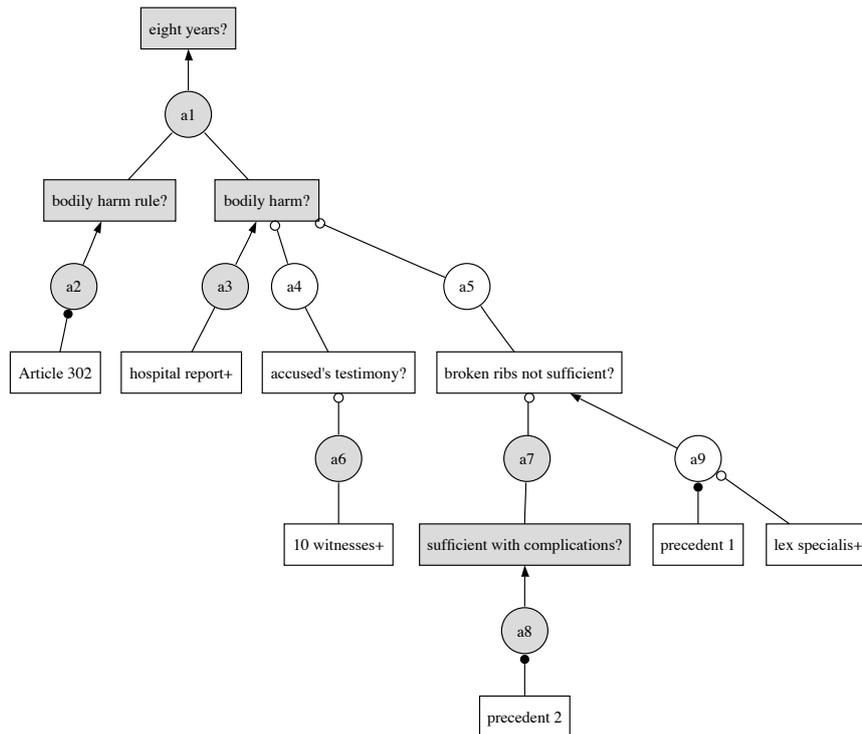


Figure 4. Reconstruction of Verheij's Grievous Bodily Harm Example

of reasoning with evidence. When completed, Carneades will support a range of argumentation use cases, including argument construction, evaluation and visualization. Although the focus of this paper was argument evaluation, it contains some hints about the direction we are heading to support argument visualization. One of our next tasks will be to refine the diagramming method used here to illustrate the argumentation framework.

REFERENCES

- [1] Katie Atkinson, Trevor Bench-Capon, and Peter McBurney, 'Arguing about cases as practical reasoning', in *Proceedings of the 10th International Conference on Artificial Intelligence and Law*, pp. 35–44, Bologna, Italy, (2005).
- [2] Tim Berners-Lee, James Hendler, and Ora Lassila, 'The semantic web', *Scientific American*, **284**(5), 34–43, (May 2001).
- [3] Floris Bex, Henry Prakken, Chris Reed, and Douglas Walton, 'Towards a formal account of reasoning with evidence: Argumentation schemes and generalizations', *Artificial Intelligence and Law*, **11**(2-3), (2003).
- [4] Jon Doyle, 'A truth maintenance system', *Artificial Intelligence*, **12**, 231–272, (1979).
- [5] Phan Minh Dung, 'On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games', *Artificial Intelligence*, **77**(2), 321–357, (1995).
- [6] Paul Edwards, *The Encyclopedia of Philosophy*, volume 1, Macmillan and Free Press, 1972.
- [7] Herbert Fiedler, 'Expert systems as a tool for drafting legal decisions', in *Logica, Informatica, Diritto*, eds., Antonio A. Martino and Fiorenza Socci Natali, 265–274, Consiglio Nazionale delle Ricerche, Florence, (1985).
- [8] Kathleen Freeman and Arthur M. Farley, 'A model of argumentation and its application to legal reasoning', *Artificial Intelligence and Law*, **4**(3-4), 163–197, (1996).
- [9] Thomas F. Gordon, 'The argument construction set — a constructive approach to legal expert systems', Technical report, German Research Institute for Mathematics and Data Processing (GMD), (1988).
- [10] Thomas F. Gordon, 'A theory construction approach to legal document assembly', in *Expert Systems in Law*, ed., Antonio A. Martino, 211–225, Amsterdam, (1992).
- [11] Thomas F. Gordon, 'The Pleadings Game — an exercise in computational dialectics', *Artificial Intelligence and Law*, **2**(4), 239–292, (1994).
- [12] Thomas F. Gordon, 'A computational model of argument for legal reasoning support systems', in *Argumentation in Artificial Intelligence and Law*, eds., Paul E. Dunne and Trevor Bench-Capon, IAAIL Workshop Series, pp. 53–64. Wolf Legal Publishers, (2005).
- [13] Thomas F. Gordon and Nikos Karacapilidis, 'The Zeno argumentation framework', in *Proceedings of the Sixth International Conference on Artificial Intelligence and Law*, 10–18, Melbourne, Australia, (1997).
- [14] Jaap Hage, 'A theory of legal reasoning and a logic to match', *Artificial Intelligence and Law*, **4**(3-4), 199–273, (1996).
- [15] Arthur C. Hastings, *A Reformulation of the Modes of Reasoning in Argumentation*, Ph.D. dissertation, Northwestern University, Evanston, Illinois, 1963.
- [16] Johan de Kleer, 'An assumption-based TMS', *Artificial Intelligence*, **28**, (1986).
- [17] Arno R. Lodder, *DiaLaw — On Legal Justification and Dialogical Model of Argumentation*, Springer, 1998.
- [18] John L. Pollock, 'How to reason defeasibly', *Artificial Intelligence*, **57**, 1–42, (1992).
- [19] Henry Prakken, 'From logic to dialectic in legal argument', in *Proceedings of the Fifth International Conference on Artificial Intelligence and Law*, 165–174, Maryland, (1995).
- [20] Horst W.J. Rittel and Melvin M. Webber, 'Dilemmas in a general theory of planning', *Policy Science*, **4**, 155–169, (1973).

- [21] Stephan E. Toulmin, *The Uses of Argument*, Cambridge University Press, 1958.
- [22] Bart Verheij, *Rules, Reasons, Arguments. Formal Studies of Argumentation and Defeat*, Ph.d., Universiteit Maastricht, 1996.
- [23] Bart Verheij, 'Dialectical argumentation with argumentation schemes: An approach to legal logic', *Artificial Intelligence and Law*, **11**(2-3), 167–195, (2003).
- [24] Bart Verheij, *Virtual Arguments*, TMC Asser Press, The Hague, 2005.
- [25] Douglas Walton, *Argumentation Methods for Artificial Intelligence in Law*, Springer, 2005.
- [26] Douglas Walton and Thomas F. Gordon, 'Critical questions in computational models of legal argument', in *Argumentation in Artificial Intelligence and Law*, eds., Paul E. Dunne and Trevor Bench-Capon, IAAIL Workshop Series, pp. 103–111, Nijmegen, The Netherlands, (2005). Wolf Legal Publishers.
- [27] Douglas Walton and Chris A. Reed, 'Argumentation schemes and enthymemes', *Synthese*, **145**, 339–370, (2005).
- [28] Douglas N. Walton, *Argument Structure : a Pragmatic Theory*, Toronto studies in philosophy, University of Toronto Press, Toronto ; Buffalo, 1996. Douglas Walton. ill. ; 24 cm.
- [29] Douglas N. Walton, *Argumentation Schemes for Presumptive Reasoning*, Erlbaum, 1996.
- [30] Douglas N. Walton, *Appeal to Expert Opinion*, Penn State Press, University Park, 1997.
- [31] Douglas N. Walton, *The New Dialectic: Conversational Contexts of Argument*, University of Toronto Press, Toronto; Buffalo, 1998. 24 cm.