

# New Experiments on Reinstatement and Gradual Acceptability of Arguments

Elfia Bezou Vrakatseli<sup>1</sup>, Henry Prakken<sup>1,2</sup>, Christian P. Janssen<sup>3</sup>

<sup>1</sup>Department of Information and Computing Sciences, Utrecht University, The Netherlands

<sup>2</sup>Faculty of Law, University of Groningen, The Netherlands

<sup>3</sup>Experimental Psychology and Helmholtz Institute, Utrecht University, The Netherlands.

e.bezouvrakatseli@students.uu.nl, h.prakken@uu.nl, c.p.janssen@uu.nl

## Abstract

This paper investigates whether empirical findings on how humans evaluate arguments in reinstatement cases support the ‘fewer attackers is better’ principle, incorporated in many current gradual notions of argument acceptability. Through three variations of an experiment, we find that (1) earlier findings that reinstated arguments are rated lower than when presented alone are replicated, (2) ratings at the reinstated stage are similar if all arguments are presented at once, compared to sequentially, and (3) ratings are overall higher if participants are provided with the relevant theory, while still instantiating imperfect reinstatement. We conclude that these findings could at best support a more specific principle ‘being unattacked is better than attacked’, but alternative explanations cannot yet be ruled out. More generally, we highlight the danger that experimenters in reasoning experiments interpret examples differently from humans. Finally, we argue that more justification is needed on why, and how, empirical findings on how humans argue can be relevant for normative models of argumentation.

## 1 Introduction

Rahwan et al. (2010) presented an empirical study of how people evaluate arguments in the context of counterarguments. Their aim was to assess how the abstract argumentation semantics of Dung (1995) treat so-called reinstatement patterns, in which an argument that is attacked by another argument is defended or ‘reinstated’ by an argument attacking the attacker, so that if there are no further arguments, the first and third argument are acceptable but the second argument must be rejected. They found that people by-and-large assess arguments according to Dung’s semantics but not fully: on a 7-point scale, the first argument was rated significantly more acceptable when presented on its own than when presented together with its attacker and defender.

There are several reasons to reconsider these experiments. A general reason is that it has been claimed that the psychological sciences face a ‘replicability crisis’ since the results of many well-known experiments appear not to be replicable (Pashler and Wagenmakers, 2012). In light of this, one aim of this paper is to test whether the results of Rahwan et al. (2010) can be replicated. A more specific reason is that since the study of Rahwan et al. appeared, the study of gradual notions of argument acceptability has become popular. These studies include probabilistic (Hunter

and Thimm, 2017), graded (Grossi and Modgil, 2019), and ranking-based (Amgoud and Ben-Naim, 2013) approaches. Some of this work refers to Rahwan et al.’s study for support of their approaches, either for gradual notions of acceptability in general (Polberg and Hunter, 2018; Hunter, Polberg, and Thimm, 2020) or for specific features of the new semantics (Grossi and Modgil, 2015, 2019; Amgoud, 2019).

In particular, Grossi and Modgil (2015) cite Rahwan et al. in support for a principle that everything else being equal, having fewer attackers is better. This principle is also a key element in several of the new semantics. For instance, all six ranking-based semantics studied by Bonzon et al. (2016) satisfy the principle of ‘void precedence’ (Amgoud and Ben-Naim, 2013), according to which an argument that has no attackers is more acceptable than an argument that has attackers, even if these attackers are counterattacked.

Accordingly, another aim of this paper is to investigate whether Rahwan et al.’s study indeed provides support for these recent developments, in particular for the ‘fewer attackers is better’ or ‘void precedence’ principle. In doing so, we will regard these formalisms not as descriptive but as prescriptive, or normative models of argumentation, that is, as modeling how people *should argue*. Our investigations are in part motivated by discussions of Cramer and Guillaume (2018a,b) and Prakken and de Winter (2018) of Rahwan et al.’s study, which give reasons to be cautious when referring to Rahwan et al. in support of the new semantics, suggesting alternative explanations for Rahwan et al.’s findings. In doing so, we do not aim to question the importance of graduality in argumentation as such. We take it for granted that graduality plays important roles in argument evaluation; the question that concerns us is how these roles can best be modelled. Moreover, we would also like to note that not all graduality semantics regard the void precedence principle as generally acceptable; for example, Bonzon et al. (2021) and Thimm and Kern-Isberner (2014) independently challenge the principle for separate reasons.

In this paper we report on three experiments in which humans evaluate arguments. The first experiment succeeded in replicating Rahwan et al.’s results on imperfect recovery from attack. The other two were aimed to test two versions of an alternative explanation for Rahwan et al.’s results suggested by Rahwan et al. and Prakken and de Winter (2018), namely, that the imperfect recovery of arguments from at-

tack is not because the participants in the experiments applied the ‘having fewer attackers is better’ principle when rating the arguments, but it is due to the specific way in which the arguments were presented to them. These experiments yielded mixed results. We evaluate the results of our experiments in light of the above-mentioned literature but also in light of the question whether empirical studies have anything to say at all about the assessment of normative theories of argumentation. Our main conclusion will be that Rahwan et al. (2010)’s results cannot (yet) be considered supporting evidence for the idea that all other things being equal, having fewer attackers is better, as embodied in the ‘void precedence’ principle, since alternative explanations for the effect they found cannot be ruled out and since a more convincing explanation is needed for why empirical findings are relevant for normative theories of argumentation.

## 2 Preliminaries

In this section the basics of Dung’s theory of abstract argumentation frameworks are summarised and applied to the reinstatement pattern that was the subject of the studies of Rahwan et al. (2010). We present Dung’s semantics in a labelling version, which is equivalent to Dung’s original semantics (Jakobovits, 2000; Caminada, 2006).

An *abstract argumentation framework* ( $AF$ ) is a pair  $\langle \mathcal{A}, \mathcal{C} \rangle$ , where  $\mathcal{A}$  is a set of *arguments* and  $\mathcal{C} \subseteq \mathcal{A} \times \mathcal{A}$  is a binary relation of *attack*. The labelling approach characterises the various semantics in terms of labellings of  $\mathcal{A}$ .

A *labelling* of an abstract argumentation framework  $\langle \mathcal{A}, \mathcal{C} \rangle$  is any assignment of either the label *in* or *out* (but not both) to zero or more arguments from  $\mathcal{A}$  such that:

1. an argument is *in* iff all arguments attacking it are *out*.
2. an argument is *out* iff it is attacked by an argument that is *in*.

Then *stable semantics* labels all arguments, while *grounded semantics* minimises and *preferred semantics* maximises the set of arguments that are labelled *in*. Relative to a semantics, an argument is *sceptically acceptable* if it is labelled *in* in all labellings, it is *rejected* if it is labelled *out* in all labellings, and it is *credulously acceptable* if it is labelled *in* in some but not all labellings.

The reinstatement pattern studied by Rahwan et al. is displayed in Figure 1. In both  $AF$ ’s argument  $A$  is sceptically

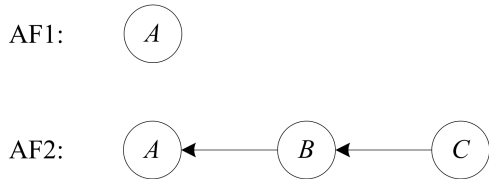


Figure 1: The reinstatement pattern

acceptable in all three semantics. With only  $A$  this is trivial since  $A$  has no attackers. When also  $B$  and  $C$  are present,  $C$

has to be made *in* by constraint (1), since it has no attackers, and  $B$  has to be made *out* by constraint (2), thus  $A$  has to be made *in* by constraint (1). Thus  $C$  reinstates  $A$  by defending  $A$  against  $B$ .

This outcome for  $AF2$  is the same if the attack from  $B$  on  $A$  is made symmetric but it changes if the attack from  $C$  on  $B$  is made symmetric (regardless whether the same is done for  $B$ ’s attack on  $A$ ). If  $C$  and  $B$  attack each other then the just-given labelling is still possible but it is not the only one: a labelling in which  $B$  is *in* and both  $A$  and  $C$  are *out* also satisfies the constraints. Both of these labellings are preferred and stable but not grounded, since the empty labelling also satisfies the constraints. Thus all three arguments are credulously acceptable in preferred and stable semantics while they are not acceptable in grounded semantics.

Rahwan et al. presented six examples to the participants in their experiments, all having the same pattern and all assumed to instantiate  $AF2$  from Figure 1. The participants were first confronted with a single argument, for instance:

$A$ : *The battery of Alex’s car is not working. Therefore, Alex’s car will halt.*

They were then asked to rate their confidence in its conclusion. Only then were they subsequently confronted with an attacker and defender, for instance:

$B$ : *The battery of Alex’s car has just been changed today. Therefore, the battery of Alex’s car is working.*

$C$ : *The garage was closed today. Therefore, the battery of Alex’s car has not been changed today.*

After both arguments, the participants were again asked to rate their confidence in the conclusion of the initial argument. After argument  $B$  their average rating of  $A$ ’s conclusion went down while after argument  $C$  was presented to them, their average rating went up again, but to a significantly lower level than after being presented with  $A$  only. Rahwan et al. concluded that their results support the notion of reinstatement but not fully, since a reinstatement argument does not fully recover from an attack.

One explanation Rahwan et al. consider for their result is in terms of an effect of ‘suspension of disbelief’, according to which participants are capable of thinking of different kinds of objections to the presented arguments but they suspend these objections for the sake of the experiment. However, when one objection is presented by the experimenter, this suspension is disrupted and some participants start to let their private beliefs ‘leak’ into their assessments of the arguments. Prakken and de Winter (2018) suggest a variation of this explanation, advocating that after being introduced to an attacker, a participant’s degree of belief in other possible attackers increases as well since the very introduction of an attacker leads them to consider other possible objections.

## 3 The Experiments

We conducted three experiments to further test these ideas. The methods of the experiments overlap and are presented together for brevity. Experiment 1 is an online replication of the study by Rahwan et al. (2010). Specifically, we test

whether rating is lower at the reinstated stage compared to the base case when arguments are presented one-by-one (cf. Rahwan et al.). Based on this replication, we then test ideas proposed by Rahwan et al. (2010) and Prakken and de Winter (2018). Specifically, experiment 2 tests whether the rating is different if all arguments (including the attack and the defense) are presented at once. Finally, experiment 3 tests what happens if first all possible scenarios are presented — i.e., generalised forms of the arguments the participants encounter during evaluation — and then the arguments are presented one-by-one. As an example of (3), the generalised form of the car battery example was

- *A car will halt if its battery is not working.*
- *A car's battery is working if it has been changed the same day.*
- *When the garage is closed, a car's battery cannot be changed.*

### 3.1 Hypotheses

We tested the following four hypotheses.

**Hypothesis 1:** When arguments are presented sequentially (experiment 1), participants' ratings for the conclusion of argument *A* in the reinstated stage are lower than in the base stage but higher than after argument *B* is presented.

The first hypothesis merely predicts a successful replication of Rahwan et al.'s results. Note that our participant number (130 aimed) is significantly higher than that used by Rahwan (20), to gain further confidence in the result.

**Hypothesis 2:** When all arguments are presented at once (experiment 2), participants' ratings for the conclusion of argument *A* are higher than the (corresponding) ratings in the reinstated stage of the first case/manner-of-presentation (where all arguments are also available but have been introduced sequentially).

The second hypothesis suggests that when all the information is presented at the same time to the participants, the confidence in the conclusion of argument *A* is higher than the corresponding confidence in the reinstated stage when arguments have been presented one-by-one. Since the introduction of an attacker may change the participant's belief in the initially presented argument even after it has been reinstated, it is possible that it is the very gradual process of presentation that influences the participant's degree of belief. To quote Rahwan et al., “[p]articipants can easily generate all sorts of objections to the arguments presented to them by the experimenter, but they suspend their disbelief in these arguments for the sake of the experiment. When one objection is presented by the experimenter herself, though, suspension of disbelief is disrupted”. Thus, if we eliminate the gradual factor of presentation, the initial suspension of disbelief may remain, since there is no stage where a *new* objection is presented that can disrupt it.

Possibly, when an attacker is introduced after one has placed their confidence in an argument, a kind of ‘breach of confidence’ is generated, one that cannot be later eradicated

(by introducing another attacker) and that has caused the disruption of the initial experiment's ‘convention/contract’ (i.e., the suspension of disbelief). Hence, if all arguments were presented at once, they could all be considered as the aforementioned ‘arguments presented by the experimenter’ and participants would suspend their disbeliefs for all of them (as suggested). Provided with all the information (i.e., all the arguments in play) at the beginning, participants can make a deliberation without the element of surprise, resulting in giving the conclusion of argument *A* a higher confidence rating than in the reinstated stage of a gradual presentation.

**Hypotheses 3a+3b:** When participants are first presented with all possible scenarios (experiment 3) — i.e., when they are presented with generalised forms of the arguments they will encounter during evaluation, before evaluating them — and are then asked to evaluate the arguments one-by-one (the same way as in experiment 1):

- a their ratings for the conclusion of argument *A* in the reinstated stage are higher than the corresponding ratings in the reinstated stage of the first experiment (where participants have not seen all the possible scenarios beforehand);
- b their ratings for the conclusion of argument *A* in the base stage are lower than the corresponding ratings in the base stage of the first experiment.

In our statistical test, we ran an Analysis Of Variance (ANOVA) with experiment (experiment 1 or 3) as between-subjects factor, and moment (base stage versus reinstated stage) as within subjects factor. Based on the hypotheses above, we would expect a significant interaction effect: rating is lower in the reinstated stage for participants in experiment 1 (compared to its base stage), whereas this is not the case for experiment 3 (i.e., no imperfect reinstatement is expected in experiment 3).

To further comment on hypotheses 3a and 3b, and extending on our thinking concerning the second hypothesis, we ought to consider another possible explanation and, thus, another manner of presentation. When a participant initially evaluates an argument, no evidence for or against its premises, inference, or conclusion has been offered, whereas after being presented with the attacker and defender, further evidence is overall provided, allowing the subject to form a more complete image of a precise situation.

Hypotheses 3a and 3b are based on Prakken and de Winter (2018), who argue that the introduction of an attacker increases the participants' degree of belief in other possible attackers, which are not explicitly ruled out in the presented arguments. They suggest that the introduction of a relevant theory prior to participants' evaluations will cause the confidence degree in the conclusion of argument *A* in the base stage to decrease (compared to ratings from the first manner of presentation) and to increase in the reinstated stage. Their suggestion is based on the assumption that if a participant was aware from the beginning of (all) the reasons argument *A* can be vulnerable, their belief in the possibility of the attacker that is presented (here, argument *B*) would increase

from the base stage, resulting in a lower rating for the conclusion *A* at that stage. By the same logic, their degree of belief in all other attackers, which are not ruled out (but neither presented) in the experiment, would have no reason to increase after the actual introduction of the attacker in the defeated stage (contrarily to when one is not initially introduced to the whole theory) and, thus, confidence in argument *A*'s conclusion would increase in the reinstated stage.

A confirmation of hypotheses 2, 3a and 3b would underline the importance of the way in which subjects are presented with arguments, proving it affects participants' confidence. Such confirmations would support the observations of Rahwan et al. and Prakken and de Winter (2018) on the possible effects of suspension of disbelief, as, then, said findings could be interpreted as a result of the two aforementioned suggested explanations and not as support for graded notions of argument acceptability.

### 3.2 Method

We conducted three experiments. In all three experiments, participants had to evaluate the acceptability status of natural language arguments, in which we followed Rahwan et al. (2010)'s method as closely as possible in terms of materials, procedure and measurement, discussed in more detail below.

**Participants** In each experiment, 130 participants took part (390 total). All were 18-65 years old. The average age was comparable between experiments (mean age for experiment 1, 2, and 3 respectively: 30, 33, and 28 years of age). All participants were volunteers, recruited through personal contact, and had no pre-knowledge on the topic of study. Participants were required to be over 18 years of age, and able to read and speak English, for which we probed participants at the start of the survey. All participants were deemed suitable according to their responses.

**Materials** The materials followed original stimuli of Rahwan et al. as close as possible. In each experiment, participants had to rate eight sets of arguments, consisting of three arguments each, where the conclusion of each next argument contradicts a premise of the preceding argument. The first six sets were taken from Rahwan et al. while the two remaining sets were added by us in a similar style. Specifically, these were:

A: *The power is out, so Claire cannot charge her phone.*

B: *The TV is playing, so the power is not out.*

C: *The TV is broken, so the TV is not playing.*

and

A: *Animals have the right to be left unharmed, so we should ban animal testing.*

B: *Animals are very dissimilar to humans, so animals do not have such a right.*

C: *Animals resemble us anatomically, physiologically, and behaviourally (e.g., recoiling from pain, fearing tormentors), therefore they are not very dissimilar to humans.*

At various points (see design), participants had to rate the acceptability of the conclusion of argument *A*. The ratings

were given on a 7-point scale ranging from *Certainly false* to *Certainly true* as in Rahwan et al. (2010).

**Design** In experiment 1, we replicate Rahwan et al. (2010). Arguments A, B, and C were added in sequence. After each added statement, participants had to rate the acceptability of the conclusion of argument *A*. Consistent with hypothesis 1 and Rahwan et al. (2010), we expect ratings to be higher after presentation of argument *A* (base stage) compared to after presentation of argument *C* (reinstated stage). This is tested with a paired t-test.

In experiment 2, all arguments are presented at once, and participants only provide one rating. We test whether this rating is different from the ratings at reinstated stage of experiment 1. Cf. hypothesis 2 we expect ratings to be higher for participants from experiment 2.

In experiment 3, for each set of arguments, participants first received a text that included generalisations of all three arguments (an example of which can be found at the beginning of section 3). They then had to rate the conclusion of argument *A* in a similar fashion as in experiment 1. As we now have a measurement at base and at reinstated stage for experiments 1 and 3, we analyze the results using an analysis of variance with experiment as between-subjects factor, and moment (base versus reinstated stage) as within-subjects factor. Conform hypothesis 3, we expect a significant interaction effect: in experiment 1 rating is lower in reinstated stage; in experiment 3 we expect there to be no or little difference between reinstated and base stage.

**Procedure** Participants did the experiment online using a Qualtrics (<https://www.qualtrics.com/>) survey. Participants were first asked a brief set of questions about their age and language capability. They then received a brief explanation of the study. Participants were then asked to rate four sets of arguments. The nature of questioning depended on which experiment they took part in (1, 2, or 3, see design). Although we had 8 sets of arguments, each participant only rated 4 sets (randomised across participants).

**Analysis** We removed data from participants whose response set was not complete (27, 34, and 20 participants in experiments 1, 2, and 3 respectively). We then calculated the average score for each rating type (reinstated stage, and base stage for experiments 1 and 3). In statistical analysis, we use alpha at .05 for significance.

### 3.3 Results

**Experiment 1 and hypothesis 1** First we test if our replication finds the same pattern of effect as Rahwan et al. (2010). A paired t-test on the data of our experiment 1 found that ratings at the base stage ( $M = 5.61, SD = 0.99, 95\% CI = [5.42, 5.81]$ ) were significantly higher compared to the reinstated stage ( $M = 5.21, SD = 0.96, 95\% CI = [5.02, 5.40]$ ),  $t(102) = 4.636, p < .001$ . Thus ratings of argument *A*'s conclusion are found to decrease after attacker *B* and increase again after counterattacker *C*, but not to the initial level, like in the original experiment of Rahwan et al. (2010). Figure 2 shows this result and also presents the values observed in Rahwan et al. (2010). It can be seen that apart from the significant difference between conditions/stages, the observed values are also com-

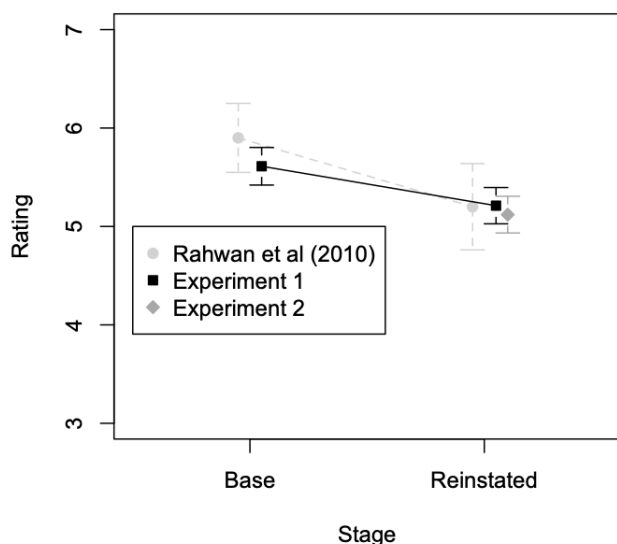


Figure 2: Rating at Base and Reinstated for 2 experiments and Rahwan et al (2010). Error bars show 95% Confidence Intervals; points are horizontally plotted slightly to the side of each other for better readability.

parable between our study and Rahwan et al. (2010) (specifically: there is a strong overlap between the error bars; the means of the two studies fall inside the region defined by the error bars). This confirms the first hypothesis, and replicates the result of Rahwan et al. (2010), this time with a considerably larger set of participants.

**Experiment 2 and hypothesis 2** Next, we test if participants give higher ratings if information is presented all at once (experiment 2) compared to sequentially (experiment 1). As the groups had unequal numbers of participants, we ran an independent Welch t-test. There was no significant effect of presentation manner on rating,  $t(196.56) = -0.683, p = .496$ . Thus, presenting all arguments at once before asking a rating of argument *A*'s conclusion does not lead to higher ratings and, so, the second hypothesis cannot be confirmed. Indeed, Figure 2 shows that the ratings in experiment 2 ( $M = 5.12, SD = 0.93$ ) are similar (i.e., means are close, error bars overlap largely).

**Experiment 3 and hypothesis 3a and 3b** Next, we test if it makes a difference if participants are provided with generalisations of all three arguments first. To this end, after checking the equality of variances of each group/experiment with Levene's test, we ran a 2 (experiment: 1 or 3) x 2 (stage: base versus reinstated) mixed ANOVA. We found a significant effect of experiment,  $F(1, 211) = 12.906, p < .001$ . There was also a significant effect of stage,  $F(1, 211) = 53.66, p < .001$ . There was no interaction between study and stage,  $F(1, 211) = 1.227, p = .269$ . Figure 3 illustrates this effect. The parallel lines suggest that in both experiment 1 and experiment 3 ratings are higher in the base stage compared to the reinstated stage, and the reduction in rating between the two is comparable (i.e.: main effect of stage).

In addition, ratings in experiment 3 were higher in general (i.e., main effect of experiment). In other words, when participants first see the possible scenarios and then rate the arguments one-by-one, they rate *A*'s conclusion higher in both the reinstated and base stage (compared to the corresponding stages of experiment 1). Thus hypothesis 3a is confirmed but hypothesis 3b is rejected. This is contrary to our expectation of an interaction effect (i.e., crossing lines in Figure 3, with the line for experiment 3 being relatively flat). The expectation was that for experiment 3 the ratings in the reinstated stage are higher than those of experiment 1 (hypothesis 3a), but that in the base stage participants from experiment 3 provided a lower rating than those in experiment 1 (hypothesis 3b). We did not observe this interaction, as hypothesis 3a was confirmed but hypothesis 3b was rejected.

## 4 Discussion

This study purported to (1) replicate the findings of Rahwan et al. (2010) and (2) investigate whether these findings support the void precedence/'fewer attackers is better' principle incorporated in many current graded notions of argument acceptability or whether alternative explanations suggested by Rahwan et al. (2010) and Prakken and de Winter (2018) undercut such support. To summarise our results, our experiment found that participants' ratings of argument *A*'s conclusion decrease after seeing attacker *B* and increase again after seeing counterattacker *C*, but not to the initial level. This confirms our hypothesis 1 and replicates Rahwan et al. (2010)'s findings. This is an important result, since replicability is one of the cornerstones of scientific method and since, as we noted in the introduction, social psychology is currently facing a replication crisis. In experiment 2 we found that presenting all arguments at once before asking a rating of argument *A*'s conclusion did not lead to higher ratings compared to those observed in the sequential study of experiment 1 (rejecting hypothesis 2). In experiment 3, we found the opposite when the participants first see the possible scenarios and then rate the arguments after seeing the arguments one-by-one (confirming hypothesis 3a). Finally, in experiment 3 we found that the participants rate *A*'s conclusion higher in the base stage as well, compared to the base stage of experiment 1 (rejecting hypothesis 3b). Thus, we did not find the interaction effect that the confirmation of both hypotheses would entail.

We now discuss various issues relevant to the question whether our results strengthen the arguments for the 'fewer attackers is better' principle.

### 4.1 Generalisation to Other Patterns

We first recall an observation of Prakken and de Winter (2018) that even if the results support a principle that 'an argument is better if it is unattacked than if it is attacked' in examples following the pattern of Figure 1, the findings cannot be used as support for the more general intuition formalised in Grossi and Modgil (2015, 2019)'s 'fewer attackers is better' principle and Amgoud and Ben-Naim (2013)'s void precedence principle, which, as noted above, is at the heart of many current gradual and ranking-based

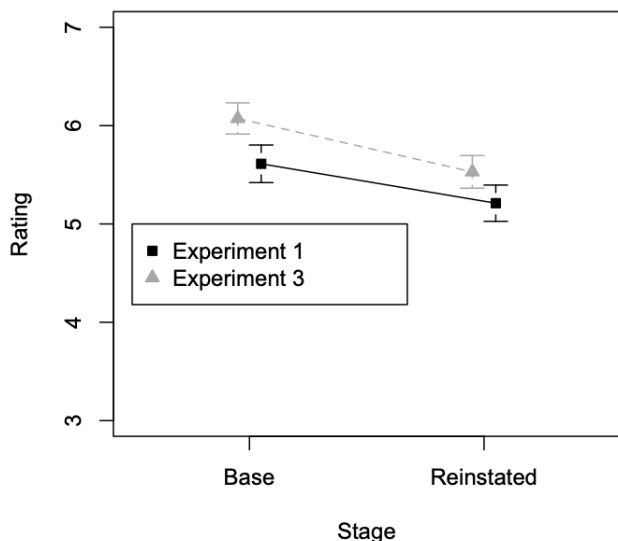


Figure 3: Rating at Base and Reinstated for experiment 1 and 3. Error bars show 95% Confidence Intervals; points are horizontally plotted slightly to the side of each other for better readability.

approaches. The point is that the more general intuition also applies to structures where, unlike in Figure 1, arguments  $A$  in  $AF1$  and  $AF2$  refer to different arguments. Neither in Rahwan et al. (2010)’s nor in our experiments examples of this kind were shown to the participants. So at best Rahwan et al. (2010)’s and our experiments confirm the special case of the void precedence/‘fewer attackers is better’ principle in which the arguments  $A$  in both  $AF$ s in Figure 1 are the same argument.

## 4.2 Suspension of Disbelief

We next note that our results cast some doubts on Rahwan et al. (2010)’s suggested explanation in terms of suspension of disbelief and its variant suggested by Prakken and de Winter (2018). Rahwan et al. do not claim that the introduction of an attacker makes the subjects think/come up with objection, but rather that it causes them to disrupt their suppressing of their already existent objections. In this study, we hypothesised that if confronted with all three arguments at the same time, participants would apply their suspension of disbelief to all the (initially) presented arguments. As our hypothesis 2 is rejected — i.e., introducing all three arguments at the same time does not have a significant effect on the subjects’ confidence in  $A$ ’s conclusion — Rahwan et al.’s explanation regarding the disruption of suspension of disbelief cannot be validated.

The same holds for Prakken and de Winter (2018)’s variant of the explanation in terms of suspension of disbelief, according to which the initial introduction of the relevant theory would have made the participants in group 3 aware of possible counterarguments from the start, unlike the participants in group 1. This should have led to the ratings for the conclusion of argument  $A$  in the base stage of group 3 being significantly lower than those of group 1, which was

our hypothesis 3b. However, this hypothesis was rejected and, surprisingly, not only are the ratings of the third group not lower in the base stage, but they are actually significantly higher. Thus, this is a case where the possibility of an attacker was present from the beginning without it influencing negatively the ratings of the argument that could be attacked. The absence of the expected interaction effect suggests that — despite the introduction of the relevant theory beforehand — the recovery was not complete in the third group either and, thus, Prakken and de Winter’s suggestion cannot explain imperfect reinstatement.

What is puzzling is the confirmation of hypothesis 3a in contrast to the rejection of hypothesis 3b, as what we expected was that the introduction of the theory would have opposite results on the base and reinstated stage. One reason why the introduction of the corresponding theory results in an increase of the ratings’ level in both stages could be that when introduced with a theory beforehand, the participant gains reassurance. Even though aware of the possibility of an attacker, when an argument is unattacked, the participant has no reason/evidence not to believe it. Thus the introduction of a *possible* attacker might in this case strengthen the attacker’s *absence in the base stage*, thus increasing confidence in the conclusion of argument  $A$ . This could even be extended to the reinstated stage: participants might feel more reassured after being presented with the instantiation of the possibilities they were originally introduced with. This could also explain why a similar effect did not appear in the second group: in the third group, a participant is originally introduced to possibilities, which are later realised, whereas in the second group a participant misses this intermediate step of reassurance. However, the results of the second group could also be explained by the task of group 2, as we will further discuss in Section 4.4.

## 4.3 Natural Language versus Formalisation

At this point, it might be thought that our findings strengthen the support for the ‘fewer attackers is better’ principle. The underlying idea here would be that the participants rated the arguments’ conclusions with this principle in mind. We first discuss whether this explanation can be accepted on the basis of Rahwan et al. (2010)’s and our experiments. Later we will discuss to which extent such empirical claims and explanations are relevant for assessing normative models of argumentation.

There is yet another alternative explanation of the results, independently suggested by Prakken and de Winter (2018) and Cramer and Guillaume (2018a,b), namely, that when rating the arguments, the participants may not have had the reinstatement pattern of Figure 1 in mind but a different pattern. All argument sets in the studies of Rahwan et al. (2010) and ourselves were such that the conclusion of argument  $B$  attacks a premise of argument  $A$  and, likewise, the conclusion of argument  $C$  attacks a premise of argument  $B$ . Consider again the car battery example from Section 2. It is not obvious that the attacks of  $B$  on  $A$  and of  $C$  on  $B$  are asymmetric: since the conclusions and premises involved in these attacks are contradictory, the attacks might also be regarded as symmetric. This is, for instance, possible in  $AS$ -

$PIC^+$  (Modgil and Prakken, 2014) in which a so-called ‘ordinary’ premise can rebut an argument with a ‘defeasible top rule’. Moreover,  $AF$ s generated in  $ASPIC^+$  include the subarguments of all arguments as separate arguments, including arguments corresponding to a premise.

Thus another plausible  $AF$  modelling of the car battery example is  $AF_2'$  as shown in Figure 4, where  $Ap$ ,  $Bp$ , and  $Cp$  are the subarguments of, respectively,  $A$ ,  $B$ , and  $C$ , consisting of their premise. Note that  $B$  and  $Ap$  attack each other since  $B$  undermines  $Ap$  (and  $A$ ) while  $Ap$  rebuts  $B$ . Likewise for the other attacks. Note also that unlike in  $AF_2$ ,

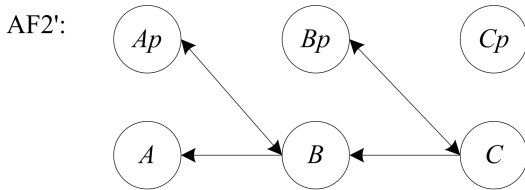


Figure 4: An alternative interpretation of Rahwan et al. (2010)’s examples

in  $AF_2'$  argument  $A$  is not skeptically acceptable. Now it is important to note that a number of participants may have interpreted the examples as in  $AF_2'$  instead of  $AF_2$ . In an experiment conducted by Cramer and Guillaume (2018a) this was indeed found to be the case. The participants who did so may have rated the conclusion of  $A$  lower in the reinstatement stage since  $A$  is only credulously acceptable in that stage.

This is an instance of a more general problem with this kind of empirical research. In experiments like these, a natural-language reasoning example is formalised, then humans are asked to express an opinion of the natural-language version of the example, and then the humans’ responses are compared to the outcome yielded by the semantics of the formalised version of the example. If there is a mismatch between the two, then it is tempting to conclude that humans do not reason according to the formal semantics but such a conclusion is premature, since the mismatch may also be caused by the fact that the formalisation does not correspond to what the humans had in mind (this point is also made by Polberg and Hunter (2018)). Formalising informal reasoning examples is far from trivial since natural language is ambiguous and the same informal way of reasoning may be formalised in the same formalism in different ways. The danger of such mismatches between a formalised example and how humans interpret its natural-language version increases the more abstract the formalism is. As noted by Prakken and de Winter (2018), directly formalising natural-language examples in abstract argumentation formalisms without being guided by a theory of the nature of arguments and their relations may result in ad-hoc modellings (or in the present case in a modelling that is not the only possible one).

This danger may also have materialised in a study of Rosenfeld and Kraus (2015), who modelled natural-language examples in a bipolar argumentation framework (an  $AF$  with also support relations) and then observed that

the participants did not assess arguments according to its semantics, including the reinstatement pattern. This result was cited by Amgoud (2019) as support for the ‘having fewer attackers is better’ principle. However, in Rosenfeld and Kraus’s examples several attack relations modelled as asymmetric can also be regarded as symmetric. For example, the arguments “We should buy an SUV; it’s the right choice for us” and “But we can’t afford an SUV, it’s too expensive” (where according to Rosenfeld and Kraus the second asymmetrically attacks the first) could by some participants be regarded as two arguments with contradictory conclusions ‘we should buy an SUV’ and ‘we should not buy an SUV’.

A related problem with such empirical reasoning experiments is that it is often hard to make the participants stick to the information that was explicitly given; often they will, either implicitly or explicitly, also take other beliefs and background information into account. Van Benthem (2008) (cited by Rahwan et al. (2010) in support of the relevance of empirical research for normative theories) notes that people in such experiments first go through a “representation” or “modelling” phase in which they construe the relevant scenario of facts and events, and only then make inferences from the construed scenario. He points at the possibility that experimenters overlook that the participants may have added information to the example in the representation phase. Other recent empirical studies in computational argumentation have also pointed at the possibly confounding effect of background information (Cerutti, Tintarev, and Oren, 2014; Polberg and Hunter, 2018; Cramer and Guillaume, 2018b, 2019).

In the present study we tried to avoid the unwanted influence of background information as follows. Overall, the arguments that were used were simple sentences and of a neutral subject matter, to avoid unwanted influence of subjective views. Moreover, the levels of confidence in the eighth set (i.e., the one regarding animal rights, which is not a neutral subject matter), do not deviate from the rest in any way. This suggests a good level of impartiality from the participants. Nevertheless, we cannot exclude the possibility that the results are partly influenced by the *content* of the arguments rather than their *relations*. In order to render such experiments less precarious, future empirical research could try to control for such issues by including manipulation checks, where separate groups of participants evaluate the arguments independently, indicate how they perceive the type and the directionality of the attacks, and so on.

#### 4.4 Order and Cognitive Load

There are further possible explanations of some of the findings. First, the results of the second group could also be explained by the task of group 2, (i.e., the version of manner of presentation that corresponded to group 2) being more challenging. As mentioned by Cramer and Guillaume (2019), a cognitively challenging task might lead to participants choosing a simplifying strategy, in this case, more likely to choose a ‘neutral’ rating (in this experiment, that would translate to a rating being closer to 4, hence being the lowest rated). The low overall ratings of argument  $A$  in group 2, along with the fact that group 2 had the highest dropout rate

(26% of the participants of the second group left the survey unfinished, compared to the 21% of the first and then 15% of the third) might be an indication that the manner of presentation in the second group was more challenging to the subjects.

Second, the imperfect recovery from attack could be a result of *order*. For example, the order of presentation may have had an effect on how the participants perceived the directionality of the attacks; it may be that attacks are more often regarded as originating from the last-presented argument. Moreover, it is possible to assume that participants' confidence in A's conclusion does not go back to its original level because the sooner we are introduced to something, the more likely we are to believe it. As observed by Polberg and Hunter (2018): "presenting a new and correct piece of information that a given person was not aware of does not necessarily lead to changing that person's beliefs". Both in our study and in Rahwan et al. (2010), the arguments were always presented in the same order. Even in group 2, where all the arguments were presented together, argument A is always first. We cannot, therefore, rule out the possibility that the order of arguments also plays a role in participants' confidence.

#### 4.5 The Jump from Is to Ought

Nevertheless, suppose that future experiments are able to reproduce Rahwan et al. (2010)'s findings in examples that unambiguously correspond to Figure 1 and in which background information has been controlled for. Then there is another hurdle to take before it can be concluded that these results support the 'having fewer attackers is better' principle as a normative principle of rational argumentation. This hurdle is that it is not immediately obvious how empirical findings on how people *actually* argue can be relevant for a normative theory on how they *should* argue. Given the growing number of empirical studies in computational argumentation (Cerutti, Tintarev, and Oren, 2014; Rosenfeld and Kraus, 2015; Cramer and Guillaume, 2018a,b, 2019; Polberg and Hunter, 2018) this question is important, but it has no simple answers.

Rahwan et al. (2010) argue that insights from psychological experiments can be relevant to the design of software agents that can argue persuasively with humans. We could think here of IBM's Debater project (Slonim, Bilu, and Alzate, 2021). They may very well be right in this: persuasiveness is a psychological phenomenon, so psychological experiments can obviously yield relevant insights on the persuasiveness of argumentation patterns. However, in our opinion formal models like Dung (1995)'s abstract argumentation theory or more concrete structured accounts like *ASPIC<sup>+</sup>*, Defeasible Logic Programming (Garcia and Simari, 2004) or assumption-based argumentation (Toni, 2014) do not aim to model *persuasiveness* of arguments. Instead they model the (nonmonotonic) *logical* status of arguments as part of a set of arguments and their logical relations of attack and support.

Rahwan et al. also argue that empirical findings on how humans actually argue are relevant for validating formal semantics of argumentation. However, they are not explicit

on when a formal semantics should be changed because of empirical findings on how humans argue and when humans should change their way of arguing to make it fit the formal semantics. One reason to change the formal semantics might be an assumption that humans by-and-large reason correctly. For example, Pollock (1986) argued that the reasoning of humans is guided by internalised rules, while Jackson (1989) argued that any descriptive attempt constitutes a "reconstruction of people's own normative ideas". However, a compelling counterexample is formed by abundant evidence that people are generally very poor at reasoning correctly with and about probabilities (Kahneman, 2011). This is generally not regarded as invalidating probability theory as a normative theory of reasoning with probabilities (here too the relevance of background information has been noted; cf. van Benthem (2008)).

One of us has argued in Prakken (2020) that there is a weaker sense in which empirical findings on how humans reason can be relevant for normative theories of reasoning. Such normative theories should not only be rationally well-founded but also 'cognitively plausible' in that it is not too difficult for people to adhere to their standards. For this reason theories of reasoning should be stated in terms that are natural to people, such as argumentation-related concepts. Such cognitively plausible normative theories may still deviate from how people actually reason, as long as they are stated in terms that are natural to people.

Applying this to the present discussion, this means that empirical research can tell us that people tend to assess arguments in gradual terms, so that it is important to develop normative theories of gradual argument evaluation. However, the specific designs of such theories cannot be based on empirical research alone but should also apply philosophical insights. In the case of gradual and ranking-based semantics, these insights must, to the best of our knowledge, still largely be developed. For instance, the only defence of the 'having fewer attackers is better' principle besides references to empirical findings that we could find is Amgoud and Ben-Naim (2013)'s claim that this principle is "natural". We suggest that a philosophical analysis should aim to clarify what is meant by argument acceptability or argument strength and should take the nature of arguments and their relations into account.

## 5 Conclusion

In this paper we returned to the experiments of Rahwan et al. (2010) on 'simple reinstatement' patterns in formal argumentation for two reasons. First, we wanted to see whether their results can be replicated. We were able to do so with a considerably larger number of participants, which is a significant result given the current concerns about replicability of results in the social sciences, specifically in social psychology. Second, we wanted to investigate with two variants of Rahwan et al.'s experiments whether empirical findings on how humans evaluate arguments in reinstatement cases can support the 'fewer attackers is better' principle incorporated in many current graded notions of argument acceptability. We can draw the following main conclusions from our investigations.



To start with, our results casted doubt on explanations suggested by Rahwan et al. (2010) and Prakken and de Winter (2018) in terms of suspension of disbelief. According to these explanations, the imperfect recovery of arguments from attack in reinstatement patterns would be due to the triggering at various moments of awareness or consideration of other counterarguments than those presented in the experiment. In our new experiments we did not find evidence for these explanations.

However, we concluded that this does not imply that the experimental results of Rahwan et al. and the present paper support the ‘fewer attackers is better’ principle. We first noted that the experiments at best support a special case of this principle, namely, ‘an argument is better if it is unattacked than if it is attacked’ (void precedence). Next we concluded that even the special case is not supported since several alternative explanations cannot yet be ruled out, such as that a number of participants may have had different attack relations in mind. More generally, we highlighted the danger that humans involved in reasoning experiments model and/or interpret examples differently than the experimenters. Finally, we argued that even if future experiments extend to the general case and can rule alternative explanations out, still convincing arguments are needed why and how empirical findings on how humans argue can be relevant for normative models of argumentation. We suggested that the importance of such empirical findings does not lie in what they say about the validity of specific reasoning patterns but in what they say about the general concepts that a normative theory should have in order to be applicable by humans. The issue concerning the jump from ‘is’ to ‘ought’ is important since the ‘having fewer attackers is better’ principle implies that it is rational for arguers to utter as many counterarguments to an argument as possible, even if these arguments are silly and can be easily refuted. Should our normative models of argumentation really encourage arguers to build their arguments on fake news and alternative facts as much as possible?

## References

- Amgoud, L., and Ben-Naim, J. 2013. Ranking-based semantics for argumentation frameworks. In Liu, W.; Subrahmanian, V.; and Wijsen, J., eds., *Scalable Uncertainty Management. SUM 2013*, number 8078 in Springer Lecture Notes in Computer Science, 134–147. Berlin: Springer Verlag.
- Amgoud, L. 2019. A replication study of semantics in argumentation. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI-19)*, 6260–6266.
- Benthem, J. v. 2008. Logic and reasoning: Do the facts matter? *Studia Logica* 88:67–84.
- Bonzon, E.; Delobelle, J.; Konieczny, S.; and Maudet, N. 2016. A comparative study of ranking-based semantics for abstract argumentation. In *Proceedings of the 30st AAAI Conference on Artificial Intelligence (AAAI 2016)*, 914–920.
- Bonzon, E.; Delobelle, J.; Konieczny, S.; and Maudet, N. 2021. A parametrized ranking-based semantics compatible with persuasion principles. *Argument and Computation* 12:49–85.
- Caminada, M. 2006. On the issue of reinstatement in argumentation. In Fischer, M.; van der Hoek, W.; Konev, B.; and Lisitsa, A., eds., *Logics in Artificial Intelligence. Proceedings of JELIA 2006*, number 4160 in Springer Lecture Notes in AI, 111–123. Berlin: Springer Verlag.
- Cerutti, F.; Tintarev, N.; and Oren, N. 2014. Formal arguments, preferences, and natural language interfaces to humans: an empirical evaluation. In *Proceedings of the 21st European Conference on Artificial Intelligence*, 207–212.
- Cramer, M., and Guillaume, M. 2018a. Directionality of attacks in natural language argumentation. In *Proceedings of the fourth Workshop on Bridging the Gap between Human and Automated Reasoning*, 40–46.
- Cramer, M., and Guillaume, M. 2018b. Empirical cognitive study on abstract argumentation semantics. In Modgil, S.; Budzynska, K.; and Lawrence, J., eds., *Computational Models of Argument. Proceedings of COMMA 2018*. Amsterdam etc: IOS Press. 413–424.
- Cramer, M., and Guillaume, M. 2019. Empirical study on human evaluation of complex argumentation frameworks. In Calimeri, F.; Leone, N.; and Manna, M., eds., *Proceedings of the 16th European Conference on Logics in Artificial Intelligence (JELIA 2019)*, number 11468 in Springer Lecture Notes in AI, 102–115. Berlin: Springer Verlag.
- Dung, P. 1995. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming, and  $n$ -person games. *Artificial Intelligence* 77:321–357.
- Garcia, A., and Simari, G. 2004. Defeasible logic programming: An argumentative approach. *Theory and Practice of Logic Programming* 4:95–138.
- Grossi, D., and Modgil, S. 2015. On the graded acceptability of arguments. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, 868–874.
- Grossi, D., and Modgil, S. 2019. On the graded acceptability of arguments in abstract and instantiated argumentation. *Artificial Intelligence* 275:138–173.
- Hunter, A., and Thimm, M. 2017. Probabilistic reasoning with abstract argumentation frameworks. *Journal of Artificial Intelligence Research* 59:565–611.
- Hunter, A.; Polberg, S.; and Thimm, M. 2020. Epistemic graphs for representing and reasoning with positive and negative influences of arguments. *Artificial Intelligence* 281:103236.
- Jackson, S. 1989. What can argumentative practice tell us about argumentation norms? In Maier, R., ed., *Norms in Argumentation. Proceedings of the Conference on Norms*, 113–122. Dordrecht/Providence RI: Foris Publication.

- Jakobovits, H. 2000. *On the Theory of Argumentation Frameworks*. Doctoral dissertation Free University Brussels.
- Kahneman, D. 2011. *Thinking, Fast and Slow*. New York: Farrar, Straus and Giroux.
- Modgil, S., and Prakken, H. 2014. The ASPIC+ framework for structured argumentation: a tutorial. *Argument and Computation* 5:31–62.
- Pashler, H., and Wagenmakers, E. 2012. Editors' introduction to the special section on replicability in psychological science: a crisis in confidence? *Perspectives on Psychological Science* 7:528–530.
- Polberg, S., and Hunter, A. 2018. Empirical evaluation of abstract argumentation: Supporting the need for bipolar and probabilistic approaches. *International Journal of Approximate Reasoning* 93:487–543.
- Pollock, J. 1986. *Contemporary Theories of Knowledge*. Littlefield, NY: Rowman & Littlefield.
- Prakken, H., and de Winter, M. 2018. Abstraction in argumentation: necessary but dangerous. In Modgil, S.; Budzynska, K.; and Lawrence, J., eds., *Computational Models of Argument. Proceedings of COMMA 2018*. Amsterdam etc: IOS Press. 85–96.
- Prakken, H. 2020. On validating theories of abstract argumentation frameworks: the case of bipolar argumentation frameworks. In *Proceedings of the 20th Workshop on Computational Models of Natural Argument*, volume 2669 of *CEUR Workshop Proceedings*, 21–30.
- Rahwan, I.; Madakkatel, M.; Bonnefon, J.-F.; Awan, R.; and Abdallah, S. 2010. Behavioural experiments for assessing the abstract semantics of reinstatement. *Cognitive Science* 34:1483–1502.
- Rosenfeld, A., and Kraus, S. 2015. Providing arguments in discussions based on the prediction of human argumentative behavior. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI 2015)*, 1320–1327.
- Slonim, N.; Bilu, Y.; and Alzate, C. 2021. An autonomous debating system. *Nature* 591:397–384.
- Thimm, M., and Kern-Isberner, G. 2014. On controversiality of arguments and stratified labelings. In Parsons, S.; Oren, N.; Reed, C.; and Cerutti, F., eds., *Computational Models of Argument. Proceedings of COMMA 2014*. Amsterdam etc: IOS Press. 413–420.
- Toni, F. 2014. A tutorial on assumption-based argumentation. *Argument and Computation* 5:89–117.