

Evaluation of a trust-modulated argumentation-based interactive decision-making tool

Elizabeth I. Sklar · Simon Parsons · Zimi Li ·
Jordan Salvit · Senni Perumal · Holly Wall ·
Jennifer Mangels

© The Author(s) 2015

Abstract The interactive *ArgTrust* application is a decision-making tool that is based on an underlying formal system of argumentation in which the evidence that influences a recommendation, or conclusion, is modulated according to values of trust that the user places in that evidence. This paper presents the design and analysis of a user study which was intended to evaluate the effectiveness of ArgTrust in a collaborative human–agent decision-making task. The results show that users’ interactions with ArgTrust helped them consider their decisions more carefully than without using the software tool.

Keywords Argumentation · Trust · Human–agent interaction

E. I. Sklar (✉) · S. Parsons
Department of Computer Science, University of Liverpool, Ashton Street, Liverpool, UK
e-mail: e.i.sklar@liverpool.ac.uk

S. Parsons
e-mail: s.d.parsons@liverpool.ac.uk

Z. Li · J. Salvit · H. Wall
Department of Computer Science, Graduate Center, City University of New York,
365 Fifth Avenue, New York, NY, USA
e-mail: zimili.sjtu@gmail.com

J. Salvit
e-mail: jordan@jordansalvit.com

H. Wall
e-mail: holly.e.wall@gmail.com

S. Perumal
Raytheon BBN Technologies, 10 Moulton Street, Cambridge, MA, USA
e-mail: senni.peri@gmail.com

J. Mangels
Department of Psychology, Baruch College, City University of New York,
55 Lexington Avenue, New York, NY, USA
e-mail: jennifer.mangels@baruch.cuny.edu

9 **1 Introduction**

10 *Argumentation* [61] is an approach to reasoning in which the focus is not just on the con-
11 clusions that are reached, but also on the data from which the conclusions are inferred, and
12 the inference steps involved. Argumentation has technical advantages over other logic-based
13 approaches to reasoning, particularly in its ability to handle inconsistent information, but
14 also seems to be a mechanism that fits well with the way that human reasoning is carried out.
15 For example, Mercier and Sperber [45] argue that the formation of reasons for and against
16 conclusions is a fundamental part of human reasoning, while Walton and Krabbe [76] cast a
17 large part of human interaction in the form of argumentation-based dialogue.

18 Argumentation has a long history in *artificial intelligence* (AI), going back at least as
19 far as 1980 [9], and has a significant, if shorter, history in *agent-based systems*. As early
20 as the 1990s [40,67], argumentation was suggested as a mechanism to extend *negotiation*
21 between agents from the simple exchange of offers to a process that allows one agent to
22 persuade another to change its position. This led to work on argumentation to underpin joint
23 planning [53],¹ and then more general approaches to argumentation-based dialogue that could
24 capture a range of dialogue types [5,56,58,59]. In work on negotiation between autonomous
25 agents, the assumption was that the use of argumentation would lead to agreements that
26 could not be reached by other means, and this was empirically verified by [57]. Subsequent
27 work has shown that argumentation also has advantages in other forms of dialogue between
28 autonomous agents [16,39]. Given the fact that argumentation appears to be a natural way for
29 humans to reason, and that fact that argumentation is beneficial in agent-agent interactions, an
30 obvious question is: *what is the effect of using argumentation in human-agent interactions?*
31 That is the question on which we focus in this paper.

32 This is not the first paper to examine matters related to this question. For example see
33 the papers in [38]. However, this paper looks at the use of argumentation in human-agent
34 interaction in a novel context: one in which the human and agent reason collaboratively,
35 using argumentation as the medium through which they represent their beliefs and thought
36 processes. The “agent” in the system described here employs formal argumentation for
37 reasoning about a scenario, and the human is given the same scenario. The agent presents its
38 conclusions and relevant evidence (i.e., its arguments) to the user, who can indicate his/her
39 level of agreement with the agent’s beliefs. In the version of the system presented here, there
40 is no dialogue about their beliefs; though future work will explore such further levels of
41 interactivity.

42 In systems of argumentation in which arguments are constructed from logical state-
43 ments [4,8,19,46], an important feature is the way in which elements of the arguments—the
44 premises and rules from which they are constructed—have a bearing on the quality of the
45 arguments. Premises may be undermined and hence defeated. Conclusions may be rebutted,
46 and rules themselves may be undercut. This relationship between the parts and the whole,
47 combined with the relationship between the trust individuals place in information and the
48 provenance of that information [21], led us to suggest the use of argumentation in situations
49 where trust in information is critical [52,55]. The key idea is that since argumentation tracks
50 the data used in deriving conclusions, if that data could be related to the sources from which
51 it comes, information about those sources could be used in reasoning about the conclusions.

52 We developed a formal argumentation system [68] that allows information about sources—
53 represented in the form of the “trust networks” that are standard in the literature of reasoning
54 about trust—to be combined with arguments. This formal system was initially implemented in

¹ Called “negotiation” in [53], but much closer to what [76] calls “deliberation”.

an inference engine called *ArgTrust* [54]. The version of *ArgTrust* evaluated here is intended as a prototype for an intelligent interactive agent that can collaborate with a user in making decisions that involve complex situations involving the analysis of data from a range of sources, not all of which can be fully trusted, and which change continuously. In the work presented here, we report on a user study designed to explore how effective *ArgTrust* is in supporting human decision-making and, since *ArgTrust* interacts with the human by providing arguments, how effective argumentation is for human-agent communication. In particular, the aim of the user study was to gather information about how people reason, how they make decisions in uncertain situations, and how they explain their decisions. Participants (i.e., human subjects) used *ArgTrust* to help them visualise a scenario and make sense of information presented that describes elements of the scenario in different ways.

The remainder of this paper is structured as follows. We start in Sect. 2 with a brief description of the *ArgTrust* system, and pointers to papers in which the reader can obtain more detail. Then, Sect. 3 describes the design of the user study, giving full details of the materials given to participants in the study. Section 4 gives the results of the study. Related work is highlighted in Sect. 5, and then Sect. 6 concludes.

2 ArgTrust

This section briefly describes *ArgTrust* and the underlying formal model.

2.1 Theoretical basis

The formal argumentation system [68] that underpins *ArgTrust* starts with the idea that we want to represent the beliefs of a set of individuals, *Ag_s*, where each $Ag_i \in Ag_s$ has access to a knowledge base, Δ_i , containing formulae in some language \mathcal{L} . An *argument* is then:

Definition 1 (*Argument*) An *argument* A from a knowledge base $\Delta_i \subseteq \mathcal{L}$ is a pair (G, p) where p is a formula of \mathcal{L} and $G \subseteq \Delta_i$ such that:

1. G is consistent;
2. $G \vdash p$; and
3. G is minimal, so there is no proper subset of G that satisfies the previous conditions.

G is called the *grounds* of A , written $G = Grounds(A)$ and p is the *conclusion* of A , written $p = Conclusion(A)$. Any $g \in G$ is called a *premise* of A . The key aspect of argumentation is the association of the grounds with the conclusion, in particular the fact that we can trace conclusions to the source of the grounds.

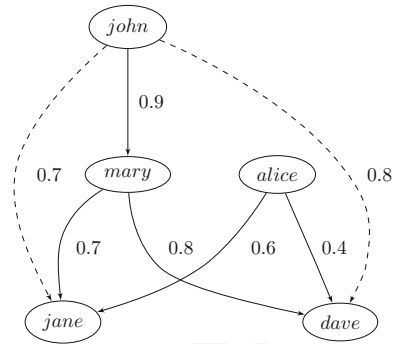
The particular language \mathcal{L} we use is the language of defeasible Horn clauses—that is, a language in which formulae are either atomic propositions p_i or formulae of the form $p_i \wedge \dots \wedge p_n \Rightarrow c$, where \Rightarrow is a defeasible rule rather than material implication. Inference in this system is by a defeasible form of generalised *modus ponens* (DGMP):

$$\frac{p_1, \dots, p_n \quad p_i \wedge \dots \wedge p_n \Rightarrow c}{c} \quad (1)$$

and if p follows from a set of formulae G using this inference rule alone, we denote this by $G \vdash p$. Given its use of defeasible Horn clauses, this argumentation system is related to that of [19].

The set of individuals, *Ag_s*, are related to each other by a social network that includes estimates of how much individual agents trust their acquaintances, as illustrated in Fig. 1.

Fig. 1 Social network. Trust is propagated using *TidalTrust* (see text)



Nodes represent individuals and links between them are annotated with the degree to which one individual trusts another, represented as values between 0 and 1. The input to the network (i.e., information known *a priori*) consists of the nodes and the solid edges. The output of the network (dashed edges) is the degree of trust inferred between any two nodes in the network. We can, for example, apply *TidalTrust* [20] to propagate trust values through the network and relate agents that are not directly connected in the social network.

In decision-making situations, argumentation can help in two ways. First, it is typical that from the data a given individual Ag_i has about a situation, we can construct a set of arguments that may conflict with each other. We might have an argument (G, p) in favour of some decision option, and another argument $(G', \neg p)$ against it (in this case, we say that the arguments *rebut* each other). We might also have a third argument $(G'', \neg g)$ where $g \in G$ is one of the grounds of the first argument (in this case we say that $(G'', \neg g)$ *undermines* (G, p)). Finally, we might have a fourth argument $(G''', \neg i)$ where i is one of the conclusions to one of the defeasible rules in (G, p) . (This is another form of rebut, rebuttal of a sub-argument.) Argumentation provides a principled way—or rather a number of alternative ways—for Ag_i to establish which of a conflicting set of arguments it is most reasonable to *accept* [7].

Second, the grounds of an argument G , can be related back to the sources of the information that constitutes the grounds. If that information comes from some individual Ag_j that Ag_i knows, then Ag_i will *believe* the information according to how much they trust Ag_j (an extension of Liau's [43] principle that you believe information from individuals that you trust). The same principle can be applied to other sources of information.² This weight can be used to resolve conflicts between arguments. It is possible to provide the decision maker with links between information that feeds into a decision and the source of that information, allowing them to explore the effect of trusting particular sources.

To see more concretely how this can be useful, let's look at a simple decision-making example, loosely based on Operation Anaconda [48] and depicted in Fig. 2. In this example, a decision is being made about whether to carry out an operation in which a combat team will move into a mountainous region to try to apprehend a high value target (HVT) believed to be in a village in the mountains.

We have the following information:

1. If there are enemy fighters in the area, then an HVT is likely to be in the area.
2. If there is an HVT in the area, and the mission will be safe, then the mission should go ahead.

² For example, military intelligence traditionally separates information into that which comes from human sources, that which comes from signals intercepts, and that which comes from imagery. All of these sources can be rated with some measure of trustworthiness.

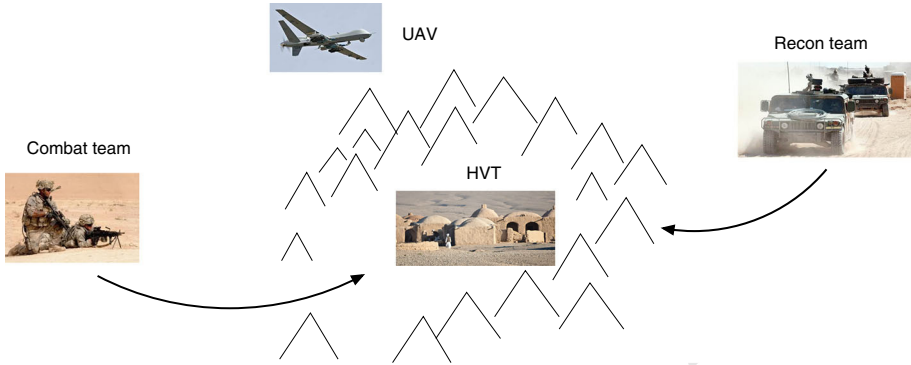


Fig. 2 The example scenario

- 129 3. If the number of enemy fighters in the area is too large, the mission will not be safe.
- 130 4. UAVs that have flown over the area have provided images that appear to show the presence
- 131 of a significant number of camp fires, indicating the presence of enemy fighters.
- 132 5. The quality of the images from the UAVs is not very good, so they are not very trusted.
- 133 6. A reconnaissance (“recon”) team that infiltrated the area saw a large number of vehicles
- 134 in the village that the HVT is thought to be inhabiting.
- 135 7. Since enemy fighters invariably use vehicles to move around, this is evidence for the
- 136 presence of many enemy fighters.
- 137 8. Informants near the combat team base claim that they have been to the area in question
- 138 and that a large number of fighters are present.
- 139 9. In addition, we have the default assumption that missions will be safe, because in the
- 140 absence of information to the contrary we believe that the combat team will be safe.

141 Thus there is evidence from UAV imaging that sufficient enemy are in the right location
 142 to suggest the presence of an HVT. There is also some evidence from informants that there
 143 are too many enemy fighters in the area for the mission to be safe.

144 We might represent this information as follows (the numbers in parentheses indicate
 145 the correspondence between the logic representations, below, and the relevant piece(s) of
 146 information, above)³:

- 147 (1) $InArea(enemy) \Rightarrow HVT$
- 148 (2) $HVT \wedge Safe(mission) \Rightarrow Proceed(mission)$
- 149 (3) $InArea(enemy) \wedge Many(enemy) \Rightarrow \neg Safe(mission)$
- 150 (4, 5) $InArea(campfires)$
- 151 (4) $InArea(campfires) \Rightarrow InArea(enemy)$
- 152 (6) $InArea(vehicles)$
- 153 (7) $InArea(vehicles) \Rightarrow Many(enemy)$
- 154 (7, 8) $Many(enemy)$
- 155 (9) $Safe(mission)$
- 156

³ While stressing that this is purely illustrative—a real model of this example would be considerably more detailed.

157 From this information, we can construct arguments such as:

$$158 \left(\left(\begin{array}{l} InArea(campfires), \\ InArea(campfires) \Rightarrow InArea(enemy), \\ InArea(enemy) \wedge Safe(mission) \Rightarrow HVT, \\ Safe(mission), \\ HVT \Rightarrow Proceed(mission) \end{array} \right), Proceed(mission) \right)$$

159 which is an argument for the mission proceeding, based on the fact that there are campfires
160 in the area, which suggest enemy fighters, that enemy fighters suggest the presence of an
161 HVT, and that the presence of an HVT (along with the default assumption that the mission will
162 be safe) suggests that the mission should go ahead. The level of belief in this argument will
163 depend on the trust in the source of the information from which the argument is constructed.
164 Since the crucial information in the argument, the presence of the campfires, is derived from
165 the UAV imaging, trust in the UAV imaging will determine the belief in the argument.

166 We can build other arguments from the available information, and, since these will conflict,
167 then compute a subset that are *acceptable*. (Approaches to this computation are discussed
168 in [7].) In this case, we can use information from the informants to build an argument that
169 there are many enemies in the area and hence the mission will not be safe:

$$170 \left(\left(\begin{array}{l} InArea(vehicles), \\ InArea(enemy), \\ InArea(vehicles) \Rightarrow Many(enemy) \\ InArea(enemy) \\ \wedge Many(enemy) \Rightarrow \neg Safe(mission) \end{array} \right), \neg Safe(mission) \right)$$

171 This conflicts with the previous argument by undermining the assumption about the mission
172 being safe. The belief in this second argument will depend, again, on the trust in the sources of
173 the information from which the argument is constructed. In this case, the crucial information
174 is that from the informants. Since information from informants is trusted less than that from
175 the UAV,⁴ the level of belief in this second argument will be less than that in the first argument.
176 Following [3], we use the degree of belief in an argument to determine which arguments are
177 *defeated*, and in this case the first argument is not defeated by the second. This, in turn means
178 that the first argument is acceptable.

179 The relation between trust in the source of an argument, defeat between arguments and
180 the computation of acceptability is explored in more detail in in [55].

181 2.2 Implementation

182 An initial version (*v1.0*) of ArgTrust was described in [69]. Here we present some aspects of
183 a more recent version, *v2.0*, which was used for the user study. Like *v1.0*, this current version
184 takes as input an XML file in a format which we sketch here. First, we have a specification of
185 how much sources of information are trusted, for example:

```
186 <trustnet>
187 <agent> recon </agent>
188 ...
189 </trust>
```

⁴ In this scenario, because informants are paid for useful information, they are widely considered to simply make up plausible information with the result that it is considered to be untrustworthy, and certainly less trustworthy than information derived from the high-resolution imaging from a UAV.

```

190 <truster> me </truster>
191 <trustee> recon </trustee>
192 <level> 0.95 </level>
193 </trust>
194 ...
195 </trustnet>

```

which specifies the individuals involved (including “me”, the decision maker) and the trust relationships between them, including the level of trust (specified as a number between 0 (no trust) and 1 (completely trustworthy)). The current implementation uses these values to compute the trust that one agent places on another using a choice of TidalTrust [20] or the mechanism described in [78].

The XML file also contains the specification of each individual’s knowledge, for example:

```

202 <beliefbase>
203 <belief>
204 <agent> recon </agent>
205 <fact> enemy_in_area </fact>
206 <level> 0.9 </level>
207 </belief>
208 ...
209 <belief>
210 <agent> me </agent>
211 <rule>
212 <premise> many_enemy </premise>
213 <conclusion> not safe </conclusion>
214 </rule>
215 <level> 1.0 </level>
216 </belief>
217 ...
218 </beliefbase>

```

Here the numbers reflect the belief each individual has in its information about the world.

From this data, and a query about a particular proposition, ArgTrust constructs arguments for that proposition by backward chaining. Once these arguments have been constructed, ArgTrust examines each formula used in the derivation of these arguments to identify if there are arguments with conclusions that attack these formulae. Each formula in those attacking arguments are then examined in turn. (And so on.) Once the full set of arguments is constructed, trust and belief are used to establish which arguments are defeated (as sketched above), and the grounded semantics [12]⁵ are applied to establish acceptability, and the conclusions labelled IN, OUT or UNDEC [7].

ArgTrust v2.0 extends the previous version [68,70] by implementing a more robust and flexible data model. ArgTrust v2.0 uses a mysql database and the Python programming language (for reasons outlined below), in place of Java (which was employed for ArgTrust v1.0). The language choice was largely made in order to simplify the recursive methods for storing the data and traversing it in different ways and because it was desirable to develop an interactive front-end that could be executed in a standard web browser. In a mysql⁶ database, we maintain arguments as a set of trees that represent the logical steps needed to arrive at the

⁵ The latest version of ArgTrust at the time of writing implements all the common semantics.

⁶ <http://www.mysql.com>

235 argument's conclusion. Thus, to return to our Operation Anaconda example, the combination
 236 of premise

$$237 \quad \textit{InArea}(\textit{campfires})$$

238 with rule

$$239 \quad \textit{InArea}(\textit{campfires}) \Rightarrow \textit{InArea}(\textit{enemy})$$

240 to infer conclusion

$$241 \quad \textit{InArea}(\textit{enemy})$$

242 would be represented as a tree in which each of the above formulae was a node, and arcs
 243 led from premise to rule to conclusion. The representation allows us to easily overlap argu-
 244 ments that share predicates or rules. Thus, if we had another argument with conclusion
 245 *InArea(enemy)*, we would represent the two arguments together as a tree with a single
 246 conclusion node.

247 Another important piece of the data model is its flexibility to receive new attributes and
 248 easily facilitate reconstructing arguments for the conclusions at hand. For example, our expe-
 249 rience is that users each have different senses of what “very trustworthy” means. Therefore,
 250 we built the system in such a way that changing values to belief levels or trust levels does not
 251 require completely reloading the scenario, instead entails just changing a parameter value.

252 The underlying ArgTrust inference engine can be invoked in four different modes: (1) as a
 253 command-line tool; (2) as a visualisation tool; (3) as an interactive decision support agent (the
 254 mode evaluated here); and (4) as a back-end reasoning engine. In command-line tool mode,
 255 a user can load an XML file, modify its contents on the ArgTrust command line, and pose
 256 queries to the inference engine. The system responds by outputting text that reports the status
 257 of arguments supporting the query. In visualisation tool mode, the system produces output in
 258 a graphical display of the resulting arguments—here the result of inference is an *argument*
 259 *graph* (see below) like that in Fig. 3. In interactive agent mode, users collaboratively step
 260 through a *decision scenario* and analyse it interactively, with help from the agent. In back-end
 261 reasoning engine mode, ArgTrust is called by another program—which might itself have an
 262 interactive front-end. Input is in the form of an XML file, as with the previous three modes;
 263 and output is also presented in the form of an XML file, where the burden of communicating
 264 the content of the output to a human user becomes the responsibility of the calling program.
 265 An example of this mode has been implemented and tested in related work involving a
 266 human–robot environment [6].

267 In visualisation and interactive agent modes, ArgTrust makes use of *argument graphs*
 268 to visualise complex scenarios and assign probabilities to all the possible outcomes. These
 269 graphs, which are distinct from the attack graphs common in the literature (also often called
 270 “argument graphs”⁷), represent the relationship between the facts and rules that make up the
 271 arguments, and the relationships between the arguments themselves. A full explanation of the
 272 graphs can be found in [68], along with the translation into graph-theoretic terms of the
 273 usual ideas of *extension* and the *acceptability* of arguments (Fig. 4).

274 The next section, below, describes the user interface developed for the interactive decision
 275 support mode. Then, Sects. 3 and 4 describe a user study designed to evaluate the effectiveness
 276 of this mode.

⁷ See, for example, Fig. 1 in [75].

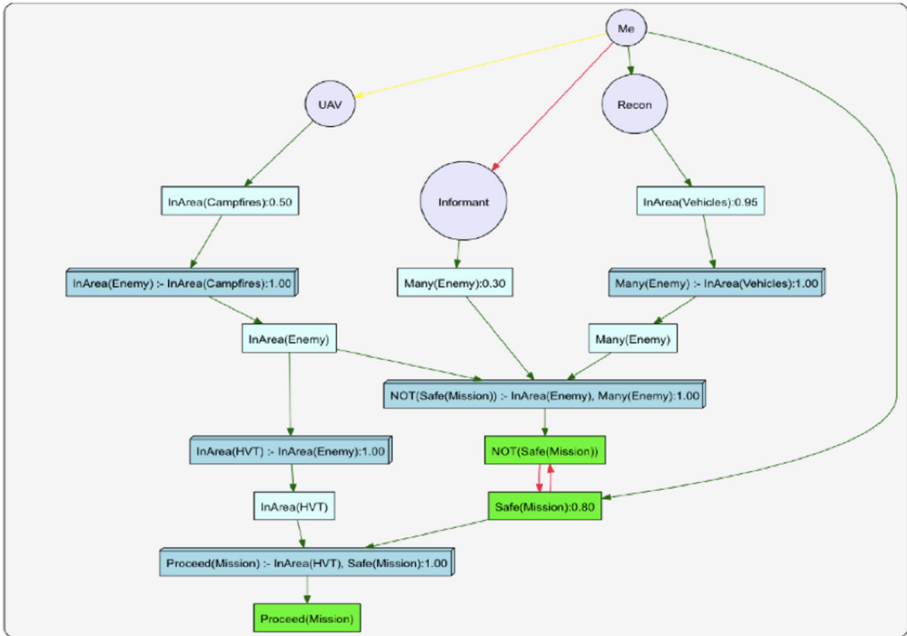


Fig. 3 ArgTrust screen: a high-level view of the argument graph

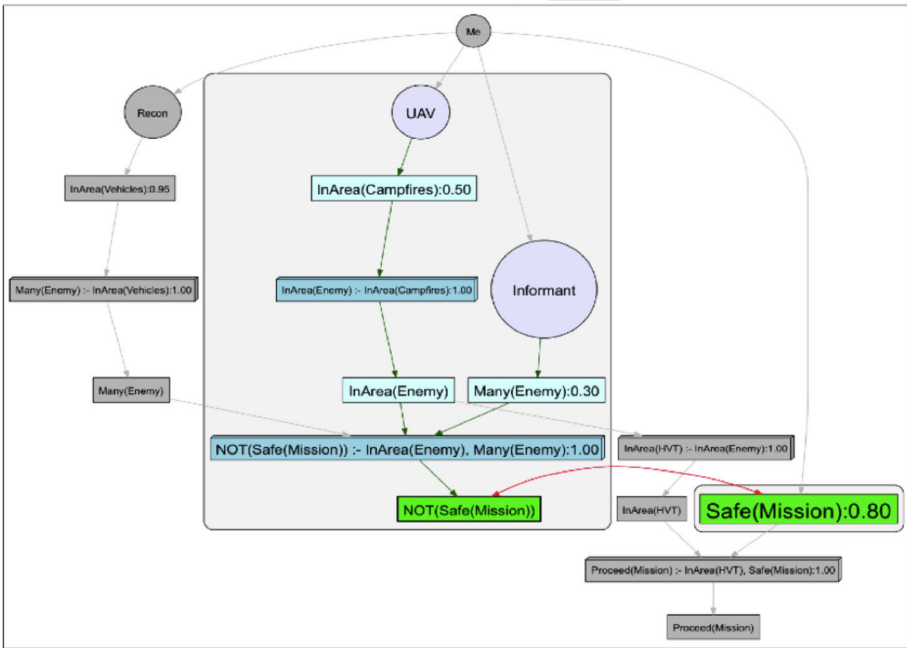


Fig. 4 ArgTrust screen: a more focussed view of one argument for the mission being unsafe

277 2.3 User interface

278 The interactive decision support mode of ArgTrust v2.0 includes an interface which allows
 279 users to manipulate the argument graphs at different levels of detail, and to focus on individual
 280 components of an argument, in order to better understand a scenario in its entirety and reach
 281 an informed conclusion. This mode was developed for the user study, and is intended to be
 282 used in a context where both the text of a narrative describing an scenario and the arguments
 283 describing a scenario are both to be presented to a test subject.

284 When used in this mode, there are four main components to the ArgTrust interface, each
 285 accessed by clicking on the following *tabs* in the user interface:

- 286 – *Review Scenario*: This component shows the text-based narrative describing a scenario
 287 and allows users to read through it before progressing. This tab is persistent, allowing
 288 users to revisit the scenario text at any time.
- 289 – *Review Trust/Beliefs/Rules*: These three tabs allow a user to change the belief level for
 290 any belief or rule, and the levels of trust in individual agents. The corresponding sentence
 291 in the scenario is displayed as well, facilitating the user's ability to set the level more
 292 accurately. After setting a belief level, the user can navigate to one of the last two tabs to
 293 see how it impacts the argument. This process can be iterative, as the user understands
 294 and learns their own process for defining belief and trust levels.
- 295 – *Trust+Beliefs*: In this part of the system, the combination of an agent's trust values
 296 and belief values are displayed to illustrate the belief value for the decision maker in a
 297 particular proposition. The goal is to fill the logical gap between setting trust values and
 298 belief values by displaying the combination of both.
- 299 – *Argument Graph*: This component displays the argument graph that corresponds to the
 300 scenario. Scenarios are broken down into arguments which are built from bits of knowl-
 301 edge (e.g., facts or evidence), rules (or beliefs), and the resulting conclusions. Facts and
 302 rules are linked to individuals (sources of information or agents). Arrows connect facts,
 303 rules and conclusions together to form a chain. A chain is referred to as an argument.
 304 (A chain can have only one arrow linking two nodes, or multiple arrows linking more
 305 than two nodes.) Each argument ends in a conclusion, and every conclusion is assigned a
 306 belief. The user can control the amount of information displayed in the graph by selecting
 307 "zoom level" and "detail" options and "focus".
 - 308 —*Zoom-level and Detail-level* controls: Located in the upper right-hand corner of the
 309 argument graph panel are the zoom and detail controls. The zoom buttons allow users to
 310 visually zoom in and out of the graph (i.e., magnifying the visual display, but not changing
 311 the content). The detail slider allows users to adjust the level of detail displayed in the
 312 graph (i.e., changing the content to be more refined or more abstract). At the highest
 313 level of detail (most abstract), only the conclusions and their corresponding beliefs are
 314 displayed. Alternatively, at the lowest level of detail (most refined), all sources, i.e.,
 315 agents, beliefs, facts, rules are shown.
 - 316 —*Focus* controls: The focus feature, located in the sidebar, enables users to focus on
 317 individual arguments of a graph. The graph updates to highlight the chosen conclusion
 318 or piece of knowledge, allowing users to focus in on that particular piece of the scenario.

319 3 User study design

320 We conducted a user study designed to evaluate the effectiveness of the ArgTrust interactive
 321 decision support mode. A primary goal of the study was to provide a preliminary assessment of

Your grandparents are coming to visit you in New York City, and they are arriving at the airport shortly. They get anxious when visiting big cities, so you promised to meet them at the airport and escort them to their hotel. You had planned to take the train to the airport straight from work. Right before you planned to leave, your co-worker tells you there was an earlier incident at a station and that train line is experiencing delays. You text a friend, who you know lives near that train line, to confirm. Your friend tells you that she left her house at the usual time and arrived to work on time, without experiencing any problems with the train. **Do you risk taking the train, which may be delayed, or do you take a taxi instead (more expensive, but quicker)?**

Fig. 5 Narrative of short, simple training scenario: “Grandparents Scenario”

the impact of ArgTrust on users’ decision-making processes. Two scenarios were developed for the user study, including narratives and logical representations of information contained in each narrative, such as in the example outlined in Sect. 2.1. One scenario is short, relatively simple and was created as a training exercise; the other scenario is longer, more complex and was built as an evaluation exercise.

The user study procedure involved multiple steps. First, participants were asked to provide demographic (e.g., gender and age) and background information (e.g., education, level of experience working with computers and decision-making tools) by filling out a Pre-survey. These data were collected for statistical purposes in order to describe the population of human subjects and to satisfy reporting requirements of funding agencies. Then participants completed a short training exercise, using the short and simple scenario mentioned above and shown in Fig. 5, to familiarise them with a formal notion of decision making under uncertainty and to give them a preliminary experience using ArgTrust.

Next, participants were presented with a text-based narrative describing a more complex scenario (the longer scenario, mentioned above and shown in Fig. 6). They were asked to analyse the details of the scenario and come to a decision about an action to take with respect to the scenario. Some of the questions required users to simply repeat information given in the scenario; this was intentional, to ensure that users had carefully read and understood the text. Once they made their decision, they were then asked to report on why and how they made that decision, via an on-line Mid-survey. Participants were asked to provide as much detail as possible regarding their thought processes.

Finally, participants were given the same scenario and asked to reconsider their decision, this time with the aid of ArgTrust. The input to ArgTrust was an XML file with contents that we extracted manually from the scenario in Fig. 6. Participants were asked to employ the user interface to interact with ArgTrust and explore the data describing the scenario, and then report on how they utilised the software in their decision-making process, by completing an on-line Post-survey. The study took approximately 60 minutes to complete.

We make a few comments on the design of the user study. First, it is likely that a *learning effect* took place between the Mid-survey and Post-survey, because participants answered questions about the same scenario in both surveys. A different study design might have had users explore two different scenarios: one for the Mid-survey (without using ArgTrust) and one for the Post-survey (after using ArgTrust). This might control better for learning effect, with respect to participants’ knowledge of the scenario (though not with respect to participants’ knowledge of the software tool). Second, a further study design might attempt to control for *order effect*: half of the participants work without ArgTrust to answer questions about the first scenario and with ArgTrust to answer questions about a second scenario; and the other half of the participants use ArgTrust to answer questions about the first scenario

A week ago, a powerful earthquake struck Brax causing widespread devastation to the country's infrastructure and leaving over 10,000 dead and over 50,000 injured. The two cities in Brax that were hit the hardest are Waga and Tapel. The Braxian Government and the UN have requested global assistance to launch the largest humanitarian relief operation ever executed. The Braxian Military, with its extensive and modern military force and airlift capability, is leading the effort and coordinating the international response. You are an Intelligence Analyst at your desk in the Operations Center of the main Forward Operating Base (FOB) in Tapel, monitoring the flow of data and reports coming in related to conditions, casualties and relief requirements. You have direct communications with the other FOB location in Waga, which was likewise affected by the earthquake. There are two rebel insurgent cells operating in the region: Reds and Lions. Each one is vying for power with the population, local and national politicians. Each one is seeking to take advantage of the situation to consolidate their political positions and establish local control with their rebel militia forces. The rebel militia forces have access to only small arms weapons and limited explosives. The rebel militias are stirring up the local population to protest the incompetence of the Braxian government. Braxian military forces are now stretched thin, trying to defend against the rebel militia forces while, at the same time, leading humanitarian rescue efforts in the wake of the earthquake. It has been 6 days since the earthquake hit Brax. Your Army Commander has asked you to answer the following Priority Intelligence Requirement (PIR): Which rebel militia cell is encouraging the most violence against the Braxian government? You have the following information (the order of the items listed is arbitrary):

- The Braxian Military reports that they have encountered many attacks/incidents of violence involving Red rebels and only some incidents of violence involving Lion rebels.
- Many incidents of violence by a rebel group imply that it is creating/encouraging much violence whereas some incidents of violence by a rebel group imply that it is not creating much violence.
- Sources of information include: Braxian Department of State, Braxian First Responders, Braxian Officials, Braxian Civilians, International Civilians and Open Media (like newspapers). Collectively, these sources of information reports only few incidents each, which makes information from them incomplete.
- Twitter feeds are inundated with reports of violence which are often contradictory. Twitter feeds are not considered very reliable.
- The Braxian Military reconnaissance reports that they have seen lots of vehicles outside the Lion Headquarters both in Tapel and in Waga.
- The presence of large number of vehicles outside a rebel militia headquarters can indicate that the rebel militia is planning many attacks/incidents of violence on relief personnel.
- Members of the rebel Lion militia who are paid by the Braxian government to inform on their comrades indicate that they have been directed to increase violence and use small arms against the Braxian military.
- A rebel group that may be planning many attacks as well as directing its members to increase violence could be a group that will create much violence.

You have to decide which rebel militia the Braxian Military efforts should focus on defending against.

Fig. 6 Narrative of longer, more complex evaluation scenario: "Humanitarian Relief Scenario"

359 and without the software for the second scenario. This might control for participants liking
 360 the software tool better because they work with it after struggling without it, or vice versa.
 361 Future studies will consider other designs such as these.

362 4 User study results

363 The user study was conducted in three sessions, where each session was conducted in a
 364 different location and involved a different set of participants. This division occurred purely due
 365 to logistics with regard to scheduling multiple sessions in which to accommodate sufficient
 366 numbers of participants. However, as mentioned below, the three groupings led to interesting
 367 distinctions with respect to the analysis of the results. **Group I** consisted of psychology
 368 undergraduate and graduate students, and the session was conducted in a university computer
 369 lab setting. **Group II** consisted of computer science undergraduate and graduate students,
 370 and the session was conducted in a university computer lab setting. **Group III** consisted of
 371 technical employees of an engineering research company, and the session was conducted
 372 in a corporate conference room. Each session followed the same procedure (outlined in
 373 Sect. 3), although the first group did not complete the Mid-survey (due to unforeseen logistical
 374 problems that occurred during the session).

375 This section discusses the results obtained by the three surveys (Pre-survey, Mid-survey
 376 and Post-survey), followed by comparative analysis across the surveys, especially the rela-
 377 tionships between answers on the Pre- and Post-surveys, and on the Mid- and Post-surveys.
 378 Questions from the Mid-survey and Post-survey are analysed in four categories: A. questions
 379 about facts (i.e., reading comprehension and paying attention to misleading questions); B.
 380 questions about applying rules found in the scenario; C. questions about trust of information
 381 sources; and D. questions about conclusions drawn from the information provided in the
 382 scenario.

383 4.1 Pre-survey

384 Twenty-two (22) participants completed the user study. Basic demographics are shown in
 385 Table 1. All participants were well-educated: 12 participants had a Masters Degree or above,
 386 and 8 had a PhD. Nobody reported previous experience with computer-based decision making
 387 tools, but almost everyone (21 of 22) claimed prior experience with data management tools,
 388 such as Microsoft Excel (20 out of 22). Only 2 of the 22 participants indicated on the
 389 Pre-survey that they had previously encountered the concepts of “logical argumentation” or
 390 “argumentation graphs”.

391 We collected the demographic data with two hypotheses in mind. The first hypothesis
 392 was that people whose education and experience was non-technical (versus technical) would
 393 find ArgTrust harder to work with. The second hypothesis was that people who were non-
 394 native (versus native) English speakers would find ArgTrust more helpful for understanding
 395 the subtleties in the narrative (which was presented in English). Table 1 tallies the number
 396 of participants in each session group who majored as *undergraduates* in Technical subjects
 397 (Computer Science, Information Technology or Engineering), as well as the number of par-
 398 ticipants who are native English speakers (at least, who speak English at home). As will
 399 be discussed below, our analysis of the results collected in the study indicates that these are
 400 relevant groupings of participants for highlighting differences in the impact and effectiveness
 401 of interacting with ArgTrust for making decisions.

402 Our first hypothesis when designing the study was that people who had no particular
 403 experience in technical subjects would find ArgTrust harder to learn how to work with,
 404 but more useful, as compared to people who had been trained in technical subjects such as
 405 Computer Science, Information Technology and/or Engineering. So we asked for participants
 406 to include information on the pre-survey about their favourite subject(s) when they were in
 407 high school and the academic subject they majored in as college undergraduates. The results

Table 1 User study participants

Group	Count	Gender		Age	
		Female	Male	18–24	25–39
Everyone	22 (100 %)	9 (41 %)	13 (59 %)	6 (27 %)	16 (73 %)
Group I	6 (27 %)	6 (100 %)	0 (0 %)	2 (33 %)	4 (67 %)
Group II	7 (32 %)	1 (14 %)	6 (86 %)	4 (57 %)	3 (43 %)
Group III	9 (41 %)	2 (22 %)	7 (78 %)	0 (0 %)	9 (100 %)
Tech	12 (55 %)	2 (17 %)	10 (83 %)	2 (17 %)	10 (83 %)
Non-Tech	10 (45 %)	7 (70 %)	3 (30 %)	4 (40 %)	6 (60 %)
English	17 (77 %)	7 (41 %)	10 (59 %)	5 (29 %)	12 (71 %)
Non-English	5 (23 %)	2 (40 %)	3 (60 %)	1 (20 %)	4 (80 %)
Group	Ethnicity				
	Asian	Black	Latino	White	
Everyone	7 (32 %)	1 (5 %)	2 (9 %)	12 (55 %)	
Group I	2 (33 %)	0 (0 %)	0 (0 %)	4 (67 %)	
Group II	3 (43 %)	1 (14 %)	1 (14 %)	2 (29 %)	
Group III	2 (22 %)	0 (0 %)	1 (11 %)	6 (67 %)	
Tech	4 (33 %)	0 (0 %)	1 (8 %)	7 (58 %)	
Non-Tech	3 (30 %)	1 (10 %)	1 (10 %)	5 (50 %)	
English	4 (24 %)	1 (6 %)	2 (12 %)	10 (59 %)	
Non-English	3 (60 %)	0 (0 %)	0 (0 %)	2 (40 %)	
Group	Academic major		Native language		
	Tech	Non-Tech	English	Non-English	
Everyone	12 (55 %)	10 (45 %)	17 (77 %)	5 (23 %)	
Group I	0 (0 %)	6 (100 %)	4 (67 %)	2 (33 %)	
Group II	3 (43 %)	4 (57 %)	5 (71 %)	2 (29 %)	
Group III	9 (100 %)	0 (0 %)	8 (89 %)	1 (11 %)	
Tech	12 (100 %)	0 (0 %)	9 (75 %)	3 (25 %)	
Non-Tech	0 (0 %)	10 (100 %)	8 (80 %)	2 (20 %)	
English	9 (53 %)	8 (47 %)	17 (100 %)	0 (0 %)	
Non-English	3 (60 %)	2 (40 %)	0 (0 %)	5 (100 %)	

There were 22 human subjects in total. The table shows participants broken down into study session, undergraduate major subject and native language groups

408 showed that 12 (55 %) of all participants in the study majored as undergraduates in Computer
 409 Science, Information Technology and/or Engineering; thus, the analysis (below) considers
 410 the participants grouped according to **Technical** majors and **Non-Technical** majors in order
 411 to evaluate this hypothesis.

412 Our second hypothesis when designing the study was that people who were non-native
 413 English speakers would find ArgTrust helpful in constructing reasoning that involved sub-
 414 tleties of the language used in the narrative (English). So we asked participants to include

Table 2 Participants' computer usage habits

Group	Several times per day	At least once per day	Several times per week	Infrequently	Never
How often do you use a computer? (i.e., laptop or desktop)					
Everyone	20 (91 %)	1 (5 %)	1 (5 %)	0 (0 %)	0 (0 %)
Group I	4 (67 %)	1 (17 %)	1 (17 %)	0 (0 %)	0 (0 %)
Group II	7 (100 %)	0 (0 %)	0 (0 %)	0 (0 %)	0 (0 %)
Group III	9 (100 %)	0 (0 %)	0 (0 %)	0 (0 %)	0 (0 %)
Tech	12 (100 %)	0 (0 %)	0 (0 %)	0 (0 %)	0 (0 %)
Non-Tech	8 (80 %)	1 (10 %)	1 (10 %)	0 (0 %)	0 (0 %)
English	15 (88 %)	1 (6 %)	1 (6 %)	0 (0 %)	0 (0 %)
Non-English	5 (100 %)	0 (0 %)	0 (0 %)	0 (0 %)	0 (0 %)
How often do you use a computer-based device? (other than a laptop or desktop)					
Everyone	20 (91 %)	1 (5 %)	1 (5 %)	0 (0 %)	0 (0 %)
Group I	6 (100 %)	0 (0 %)	0 (0 %)	0 (0 %)	0 (0 %)
Group II	6 (86 %)	1 (14 %)	0 (0 %)	0 (0 %)	0 (0 %)
Group III	8 (89 %)	0 (0 %)	1 (11 %)	0 (0 %)	0 (0 %)
Tech	11 (92 %)	0 (0 %)	1 (8 %)	0 (0 %)	0 (0 %)
Non-Tech	9 (90 %)	1 (10 %)	0 (0 %)	0 (0 %)	0 (0 %)
English	15 (88 %)	1 (6 %)	1 (6 %)	0 (0 %)	0 (0 %)
Non-English	5 (100 %)	0 (0 %)	0 (0 %)	0 (0 %)	0 (0 %)
How often do you use social media? (e.g., Facebook, Twitter)					
Everyone	6 (27 %)	2 (9 %)	6 (27 %)	7 (32 %)	1 (5 %)
Group I	2 (33 %)	1 (17 %)	2 (33 %)	0 (0 %)	1 (17 %)
Group II	3 (43 %)	1 (14 %)	1 (14 %)	2 (29 %)	0 (0 %)
Group III	1 (11 %)	0 (0 %)	3 (33 %)	5 (56 %)	0 (0 %)
Tech	2 (17 %)	1 (8 %)	3 (25 %)	6 (50 %)	0 (0 %)
Non-Tech	4 (40 %)	1 (10 %)	3 (30 %)	1 (10 %)	1 (10 %)
English	5 (29 %)	1 (6 %)	5 (29 %)	5 (29 %)	1 (6 %)
Non-English	1 (20 %)	1 (20 %)	1 (20 %)	2 (40 %)	0 (0 %)

information on the pre-survey about the language that they speak at home (“What language(s) do you and your parents speak at home?”). Although 17 (77 %) of participants speak English at home, many of the participants also speak another language; in fact, 15 (68 %) of participants speak languages other than English at home. Overall, 7 (32 %) people indicated that they live in bi-lingual households, and one person indicated that they speak four languages at home. The analysis (below) considers the participants grouped according to **English** speaking and **Non-English** speaking homes in order to evaluate this hypothesis.

We also asked participants about computer usage in the Pre-survey, in order to whether their level of familiarity with technology. The results are shown in Table 2. These indicate that all members of the cohort involved in the study were quite familiar with computers and technology in general, regardless of whether they had studied a technical subject as an undergraduate and regardless of their native language. Use of social media was more varied across the cohort. As a result, the analysis that follows does not attempt to derive any

428 correlations between computer usage habits and results with respect to understanding of the
429 test scenario or ease-of-use interacting with ArgTrust.

430 4.2 Mid-survey

431 The Mid-survey was designed to reflect users' understanding of the scenario after only
432 reading the text narrative (i.e., before interacting with ArgTrust). Some sample questions
433 are discussed here. Note that the participants from Group I did not complete the Mid-survey,
434 due to logistical issues when the study was administered at that site, so only 16 participants
435 (Groups II and III) completed the Mid-survey. Percentages reported in this section are thus
436 computed out of 16 instead of 22 users.

437 Users responded to six multiple choice questions (with possible answers of TRUE, FALSE
438 and Inconclusive), followed by an indication of their confidence in their answer, ranging
439 from 1 (least confident) to 10 (most confident). Users also responded to questions about their
440 trust in the informants depicted in the scenario and the likelihood that an intermediate con-
441 clusion drawn from some evidence provided in the scenario is valid. Finally, users responded
442 to the ultimate question posed in the scenario narrative:

443 You have to decide which rebel militia the Braxian Military efforts should focus
444 on defending against.

445 along with their confidence in their answer.

446 4.2.1 Mid-survey questions about facts

447 The multiple choice questions were designed to determine how carefully the users read
448 the narrative text and how well they understood the scenario. Specifically, we are looking
449 for whether users extracted facts (or evidence, to use the argumentation terminology) and
450 implications associated with the facts (or rules, again using the argumentation terminology).
451 Table 3 tallies the responses to Mid-survey questions about facts presented in the scenario.

452 In answer to the question:

453 The Red rebel militia is operating in the area devastated by an earthquake.

454 none of the users answered FALSE. Only just over half felt confident with a TRUE answer
455 (63 %) and the rest were unable to reach a conclusion (38 %). Close examination of the
456 narrative produces two sentences that clearly relate to this question:

457 A week ago, a powerful earthquake struck Brax...

458 There are two rebel insurgent cells operating in the region: Reds and Lions.

459 But there is indeed ambiguity with respect to the phrase “in the region”—does “region”
460 refer to the same area where the earthquake has occurred?

461 In contrast, in response to question:

462 Many vehicles have been seen near Lion headquarters in Waga.

463 all the users answered TRUE, and their confidence was high (9.06). Referring back to the
464 narrative, one of the bulleted items states:

465 – The Braxian Military reconnaissance reports that they have seen lots of
466 vehicles outside the Lion Headquarters both in Tapel and in Waga.

Table 3 Responses to Mid-survey questions about facts

Group	TRUE	FALSE	Inconclusive	Confidence
The Red rebel militia is operating in the area devastated by an earthquake				
Everyone	10 (63 %)	0 (0 %)	6 (38 %)	8.12 (1.63)
Group II	5 (71 %)	0 (0 %)	2 (29 %)	7.57
Group III	5 (56 %)	0 (0 %)	4 (44 %)	8.56
Tech	7 (58 %)	0 (0 %)	5 (42 %)	8.08
Non-Tech	3 (75 %)	0 (0 %)	1 (25 %)	8.25
English	8 (62 %)	0 (0 %)	5 (38 %)	8.31
Non-English	2 (67 %)	0 (0 %)	1 (33 %)	7.33
Many vehicles have been seen near Lion headquarters in Waga				
Everyone	16 (100 %)	0 (0 %)	0 (0 %)	9.06 (0.93)
Group II	7 (100 %)	0 (0 %)	0 (0 %)	9.57
Group III	9 (100 %)	0 (0 %)	0 (0 %)	8.67
Tech	12 (100 %)	0 (0 %)	0 (0 %)	8.83
Non-Tech	4 (100 %)	0 (0 %)	0 (0 %)	9.75
English	13 (100 %)	0 (0 %)	0 (0 %)	9.00
Non-English	3 (100 %)	0 (0 %)	0 (0 %)	9.33
Twitter feeds are unreliable sources of information				
Everyone	12 (75 %)	1 (6 %)	3 (19 %)	8.12 (1.78)
Group II	5 (71 %)	1 (14 %)	1 (14 %)	8.86
Group III	7 (78 %)	0 (0 %)	2 (22 %)	7.56
Tech	9 (75 %)	0 (0 %)	3 (25 %)	8.08
Non-Tech	3 (75 %)	1 (25 %)	0 (0 %)	8.25
English	9 (69 %)	1 (8 %)	3 (23 %)	8.00
Non-English	3 (100 %)	0 (0 %)	0 (0 %)	8.67
Vehicles belonging to Tiger rebels have been seen near the Brax Military headquarters				
Everyone	1 (6 %)	10 (63 %)	5 (31 %)	9.12 (1.20)
Group II	1 (14 %)	4 (57 %)	2 (29 %)	8.86
Group III	0 (0 %)	6 (67 %)	3 (33 %)	9.33
Tech	0 (0 %)	8 (67 %)	4 (33 %)	9.25
Non-Tech	1 (25 %)	2 (50 %)	1 (25 %)	8.75
English	1 (8 %)	7 (54 %)	5 (38 %)	9.00
Non-English	0 (0 %)	3 (100 %)	0 (0 %)	9.67
The Lions have access to chemical weapons				
Everyone	0 (0 %)	7 (44 %)	9 (56 %)	9.00 (1.10)
Group II	0 (0 %)	5 (71 %)	2 (29 %)	8.86
Group III	0 (0 %)	2 (22 %)	7 (78 %)	9.11
Tech	0 (0 %)	5 (42 %)	7 (58 %)	9.08
Non-Tech	0 (0 %)	2 (50 %)	2 (50 %)	8.75
English	0 (0 %)	5 (38 %)	8 (62 %)	9.00
Non-English	0 (0 %)	2 (67 %)	1 (33 %)	9.00

Standard deviation is shown for averages across all participants (in parentheses)

467 The unanimous consensus could stem from two reasons: first, readers tend to pay more
 468 attention to details that are enumerated in lists; and second, users trusted the Braxian Military
 469 as the source of this information more than other sources (such as informants).

470 The first reason (readers pay attention to details presented in lists) is backed up by the
 471 results shown in response to the question:

472 Twitter feeds are unreliable sources of information.

473 The following excerpt from the narrative provides evidence for the answer:

474 – Twitter feeds are inundated with reports of violence which are often contra-
 475 dictory. Twitter feeds are not considered very reliable.

476 Here, most users provided a TRUE answer (75 %), though a few could not draw a conclusion.
 477 One participant answered FALSE, which we could interpret as an indication that that person
 478 did not read the narrative carefully or fully comprehend the details; however, this person also
 479 indicated in the Pre-survey that they used social media several times a day, which may bias
 480 his/her answer. Note that there were other participants who also indicated that they use social
 481 media several times a day but still answered TRUE to this question about the unreliability
 482 of Twitter.

483 Some questions were designed to be misleading with respect to the facts presented in the
 484 scenario. Consider the statement from the Mid-survey:

485 Vehicles belonging to the Tiger rebel militia have been seen near the Braxian
 486 Military headquarters.

487 This statement indicates a group called “Tiger”, which does not exist in the scenario and
 488 was included as a contrast to the previous question which asks about “Lion”—attempting to
 489 highlight whether users confused the names “Lion” and “Tiger”. More than half the users
 490 entered FALSE (63 %), indicating that most were paying attention, one user entered TRUE,
 491 indicating that they were not paying attention; as well, a non-trivial number of users (5)
 492 thought the information was Inconclusive. This could either be because they didn’t read the
 493 scenario carefully, or because they thought that there might really be another rebel group
 494 called Tiger, but that the scenario did not provide any conclusive evidence about the Tiger
 495 group.

496 The statement:

497 The Lions have access to chemical weapons.

498 is in clear contraction to the following statement in the narrative:

499 The rebel militia forces have access to only small arms weapons and limited
 500 explosives.

501 Nonetheless, less than half the users (44 %) answered FALSE, while the remainder could
 502 not draw any conclusion (56 %); and confidence overall was high (9.00).

503 4.2.2 Mid-survey questions about applying rules

504 One question attempts to measure how well the users applied rules presented in the scenario:

505 How likely is it that vehicles outside the Lion militia HQ indicate that the group
 506 is planning an attack?

507 which puts together this fact:

Table 4 Responses to Mid-survey questions about applying rules, ranging from 1 (not likely at all) to 10 (very likely)

Group	Average likelihood (1–10)
How likely is it that vehicles outside the Lion militia headquarters indicate that the group is planning an attack?	
Everyone	6.75 (1.61)
Group II	7.00
Group III	6.56
Tech	7.00
Non-Tech	6.00
English	6.31
Non-English	8.67

Table 5 Responses to Mid-survey questions about trust, ranging from 1 (not much trust) to 10 (very much)

Group	Average level of trust (1–10)
How much do you trust the paid Lion militia informants?	
Everyone	5.69 (1.66)
Group II	6.00
Group III	5.44
Tech	5.58
Non-Tech	6.00
English	5.85
Non-English	5.00

508 The Braxian Military reconnaissance reports that they have seen lots of vehicles
509 outside the Lion Headquarters...

510 and this rule:

511 The presence of large number of vehicles outside a rebel militia headquarters
512 can indicate that the rebel militia is planning many attacks...

513 from the narrative. The results (see Table 4) show that users did not assimilate this information
514 well from reading the scenario, because the mean answer (and standard deviation) were
515 6.75 (1.61), on a scale of 1 (not likely at all) to 10 (very likely).

516 *4.2.3 Mid-survey questions about trust*

517 One question attempts to measure how much users *trust* sources of information:

518 How much do you trust the paid Lion militia informants?

519 The mean answer (and standard deviation) were 5.69 (1.66), on a scale of 1 (not much at
520 all) to 10 (very much). Results are tallied in Table 5.

521 *4.2.4 Mid-survey questions about conclusions*

522 Finally, one question asks about users’ conclusions with respect to the overall predicament
523 posed by the narrative:

Table 6 Responses to Mid-survey questions about conclusion drawn

Group	Lions	Reds	Confidence
Which rebel militia is the most violent?			
Everyone	7 (44 %)	9 (56 %)	7.19 (1.17)
Group II	4 (57 %)	3 (43 %)	7.71
Group III	3 (33 %)	6 (67 %)	6.78
Tech	5 (42 %)	7 (58 %)	6.83
Non-Tech	2 (50 %)	2 (50 %)	8.25
English	5 (38 %)	8 (62 %)	7.31
Non-English	2 (67 %)	1 (33 %)	6.67

524 Which rebel militia is the most violent, implying that Braxian Military efforts
 525 should focus on defending against attacks from that insurgent group (as
 526 opposed to other insurgents)?

527 The responses are tallied in Table 6. Results are fairly evenly split (44 : 56) between the
 528 two groups (*Lions*) and (*Reds*), and confidence is moderate (7.19)—lower with respect to
 529 questions about facts (which averaged 8.7 confidence measures).

530 4.3 Post-survey

531 The Post-survey is divided into two parts. The first part contains a series of questions that,
 532 like the Mid-survey, are designed to reflect how well the users understand the scenario
 533 and are supported in making decisions with the information presented. The second part
 534 contains a number of questions that are designed to assess the users' impressions of ArgTrust,
 535 independent from the particular scenario. Here, we discuss both parts. Comparative analysis
 536 of answers given in the Post-survey and Mid-survey is deferred to Sect. 4.3.6. We begin, next,
 537 by presenting the Post-survey responses.

538 4.3.1 Post-survey, Part 1

539 The first part of the Post-survey was designed to reflect users' understanding of the scenario
 540 after interacting with ArgTrust. Some sample questions are discussed here. Percentages
 541 reported in this section are computed out of 22 users, since all users completed the Post-
 542 survey.

543 Similarly to the Mid-survey, users responded to four multiple choice questions (with possi-
 544 ble answers of TRUE, FALSE and Inconclusive), as well as indicating their confidence
 545 in their answer, ranging from 1 (least confident) to 10 (most confident). Users also responded
 546 to questions about their *trust* in the informants depicted in the scenario and the *likelihood*
 547 that an intermediate conclusion drawn from some evidence provided in the scenario is valid.
 548 Finally, users responded to the ultimate question posed in the scenario narrative:

549 You have to decide which rebel militia the Braxian Military efforts should focus
 550 on defending against.

551 along with their confidence in their answer.

552 4.3.2 *Post-survey questions about facts*

553 As with the Mid-survey, the multiple choice questions in the Post-survey were designed
 554 to determine how carefully the users read the narrative text and how well they understood
 555 the scenario, in particular after interacting with ArgTrust to work with the facts and rules
 556 contained in the narrative. The responses are tallied in Table 7.

557 The following statement was contained in the Post-survey:

558 Information from the Braxian Military reconnaissance team indicates that there
 559 are many rebel militia forces in the areas devastated by the earthquake.

560 Evidence supporting this observation is contained in the scenario narrative:

561 A week ago, a powerful earthquake struck Brax...

562 There are two rebel insurgent cells operating in the region: Reds and Lions.

563 which can reasonably derive a TRUE response to the question. As Table 7 indicates, more
 564 than half of participants (64 %) responded with TRUE. However, most of the rest (7 people,
 565 32 %) were unable to draw a conclusion.

566 Another Post-survey question:

567 Vehicles being seen near a rebel headquarters implies that rebels are planning
 568 an attack.

569 paraphrases an item from the narrative:

570 – The presence of large number of vehicles outside a rebel militia headquar-
 571 ters can indicate that the rebel militia is planning many attacks/incidents of
 572 violence on relief personnel.

573 However, a literal interpretation of the narrative phrase “can indicate” is less conclusive
 574 than the phrasing in the question: “implies”. Perhaps this is why more users responded to the
 575 question with an Inconclusive answer, as shown in Table 7. Here we note that all hesitation
 576 with respect to this fact was reported by those who speak English at home.

577 Confidence with respect to questions about facts was lower than the Mid-survey, 7.78 on
 578 average, versus 8.7 average confidence on the Mid-survey.

579 4.3.3 *Post-survey questions about applying rules*

580 As with the Mid-survey, one question attempts to measure how well the users applied rules
 581 presented in the scenario:

582 How likely is it that vehicles outside the Lion militia HQ indicate that the group
 583 is planning an attack?

584 Results are shown in Table 8. The mean answer (and standard deviation) were 6.91 (1.72),
 585 on a scale of 1 (not likely at all) to 10 (very likely). This is an increase from the Mid-survey
 586 (which was 6.75), though not statistically significant.

587 4.3.4 *Post-survey questions about trust*

588 Again, as with the Mid-survey, one question attempts to measure how much users *trust*
 589 sources of information:

Table 7 Responses to Post-survey questions about facts

Group	TRUE	FALSE	Inconclusive	Confidence
Information from the Braxian Military reconnaissance team indicates that there are many rebel militia forces in the areas devastated by the earthquake				
Everyone	14 (64 %)	1 (5 %)	7 (32 %)	7.64 (1.99)
Group I	2 (33 %)	1 (17 %)	3 (50 %)	6.50
Group II	6 (86 %)	0 (0 %)	1 (14 %)	8.14
Group III	6 (67 %)	0 (0 %)	3 (33 %)	8.00
Tech	9 (75 %)	0 (0 %)	3 (25 %)	8.25
Non-Tech	5 (50 %)	1 (10 %)	4 (40 %)	6.90
English	11 (65 %)	0 (0 %)	6 (35 %)	7.88
Non-English	3 (60 %)	1 (20 %)	1 (20 %)	6.80
Vehicles being seen near a rebel headquarters implies that rebels are planning an attack				
Everyone	10 (45 %)	1 (5 %)	11 (50 %)	7.91 (1.57)
Group I	3 (50 %)	1 (17 %)	2 (33 %)	7.67
Group II	3 (43 %)	0 (0 %)	4 (57 %)	8.57
Group III	4 (44 %)	0 (0 %)	5 (56 %)	7.56
Tech	6 (50 %)	0 (0 %)	6 (50 %)	7.92
Non-Tech	4 (40 %)	1 (10 %)	5 (50 %)	7.90
English	5 (29 %)	1 (6 %)	11 (65 %)	8.00
Non-English	5 (100 %)	0 (0 %)	0 (0 %)	7.60

Standard deviation is shown for averages across all participants (in parentheses)

Table 8 Responses to Post-survey questions about applying rules, ranging from 1 (not likely at all) to 10 (very likely)

Group	Average likelihood (1–10)
How likely is it that vehicles outside the Lion militia headquarters indicate that the group is planning an attack?	
Everyone	6.91 (1.72)
Group I	7.50
Group II	6.86
Group III	6.56
Tech	6.83
Non-Tech	7.00
English	6.47
Non-English	8.40

590 How much do you trust the paid Lion militia informants?

591 Results are tallied in Table 9. The mean answer (and standard deviation) were 5.36 (1.71), on
 592 a scale of 1 (not much at all) to 10 (very much). This is a small decrease from the Mid-survey
 593 (5.69), though not statistically significant. It is notable, however, that all participants *except*
 594 those from non-English speaking households lowered their level of trust between the Mid-
 595 and Post-surveys.

Table 9 Responses to Post-survey questions about trust, ranging from 1 (not much trust) to 10 (very much)

Group	Average level of trust (1–10)
How much do you trust the paid Lion militia informants?	
Everyone	5.36 (1.71)
Group I	5.33
Group II	5.71
Group III	5.11
Tech	5.00
Non-Tech	5.80
English	5.29
Non-English	5.60

596 4.3.5 Post-survey questions about conclusions

597 Two of the multiple-choice questions in Part 1 of the Post-survey ask users to draw their own
598 intermediate conclusions based on the narrative:

599 Red leaders have tweeted that they are planning an attack.

600 and

601 Braxian Military should prioritize rescue operations over attacks on rebel militia.

602 The answers, tallied in Table 10 indicate that users were unable to draw their own intermediate
603 conclusions. This is understandable, because relatively little information pertinent to either
604 question was provided in the narrative text, and providing a TRUE or FALSE answer to
605 either question would require the user to make unsubstantiated inferences. This is a positive
606 effect of interacting with ArgTrust, because users are given the opportunity to modulate the
607 facts and rules contained in the narrative according to their own beliefs and trust values. The
608 act of working directly with the content forces users to be careful not to draw conclusions
609 for which there is no supporting evidence in the scenario.

610 Finally, users are asked to draw a conclusion with respect to the overall predicament posed
611 by the narrative:

612 You have to decide which rebel militia the Braxian Military efforts should focus
613 on defending against.

614 The results are displayed in Table 11. The responses lean heavily toward defending against the
615 Reds militia (82 %), though confidence is moderate (6.55). It is notable that confidence was
616 lower for *all* groupings of participants and more sharply lowered by non-technical participants
617 and participants from non-English speaking homes.

618 4.3.6 Changes from mid-survey to post-survey

619 This section highlights responses that changed between the Mid-survey and the Post-survey.
620 Note that because (as mentioned earlier) only participants in Groups II and III completed the
621 Mid-survey, percentages presented in this section are computed out of 16 instead of 22.

622 We report on three types of changes. The first type is the extent to which users changed
623 their opinions about facts and/or rules presented in the scenario. The second type is the extent

Table 10 Responses to post-survey questions about intermediate conclusions

Group	TRUE	FALSE	Inconclusive	Confidence
Red leaders have tweeted that they are planning an attack				
Everyone	0 (0 %)	8 (36 %)	14 (64 %)	8.50 (2.44)
Group I	0 (0 %)	2 (33 %)	4 (67 %)	8.00
Group II	0 (0 %)	4 (57 %)	3 (43 %)	8.57
Group III	0 (0 %)	2 (22 %)	7 (78 %)	8.78
Tech	0 (0 %)	4 (33 %)	8 (67 %)	8.92
Non-Tech	0 (0 %)	4 (40 %)	6 (60 %)	8.00
English	0 (0 %)	4 (24 %)	13 (76 %)	8.71
Non-English	0 (0 %)	4 (80 %)	1 (20 %)	7.80
Braxian Military should prioritize rescue operations over attacks on rebel militia				
Everyone	5 (23 %)	2 (9 %)	15 (68 %)	7.45 (1.92)
Group I	1 (17 %)	0 (0 %)	5 (83 %)	7.50
Group II	4 (57 %)	1 (14 %)	2 (29 %)	8.29
Group III	0 (0 %)	1 (11 %)	8 (89 %)	6.78
Tech	2 (17 %)	1 (8 %)	9 (75 %)	7.33
Non-Tech	3 (30 %)	1 (10 %)	6 (60 %)	7.60
English	3 (18 %)	2 (12 %)	12 (71 %)	7.41
Non-English	2 (40 %)	0 (0 %)	3 (60 %)	7.60

Table 11 Responses to Post-survey questions about conclusion drawn

Group	Lions	Reds	Confidence
Which rebel group should the Braxian Military provide defense against?			
Everyone	4 (18 %)	18 (82 %)	6.55 (1.63)
Group I	2 (33 %)	4 (67 %)	5.67
Group II	1 (14 %)	6 (86 %)	7.14
Group III	1 (11 %)	8 (89 %)	6.67
Tech	1 (8 %)	11 (92 %)	6.67
Non-Tech	3 (30 %)	7 (70 %)	6.40
English	2 (12 %)	15 (88 %)	6.94
Non-English	2 (40 %)	3 (60 %)	5.20

to which users changed their confidence in their answers. The third type is the extent to which users changed the number of **Inconclusive** answers to conclusive answers (i.e., **TRUE** or **FALSE**).

With respect to reflecting the facts and/or rules presented in the scenario, the following changes in opinion were detected:

- In response to questions about whether there are rebel groups in the same region as the earthquake, two (13 %) users changed their opinions from **FALSE** to **TRUE**.
- Half the users (50 %) also changed their opinion about how likely is it that vehicles outside the Lion militia HQ indicate that the group is planning an attack. The mean (and standard deviation) rose from 6.75 (1.61) to 6.69 (1.62).

Table 12 Differences in confidence levels, from mid-survey to post-survey

Group	Mid-survey		Post-survey		Change
Everyone	8.48	(0.70)	7.81	(0.90)	-0.67
Group II	8.57	(0.70)	8.14	(0.78)	-0.43
Group III	8.41	(0.72)	7.56	(0.94)	-0.86
Tech	8.45	(0.73)	7.82	(0.94)	-0.64
Non-Tech	8.57	(0.68)	7.80	(0.91)	-0.77
English	8.45	(0.70)	7.72	(0.91)	-0.73
Non-English	8.62	(0.81)	8.20	(0.87)	-0.42

Negative *change* means that level of confidence decreased. Only data from participants who completed both mid-survey and post-survey is included

Table 13 Percentage of multiple choice questions that were answered Inconclusive

Group	Mid-survey		Post-survey		Change
Everyone	0.34	(0.20)	0.52	(0.28)	0.17
Group II	0.10	(0.18)	0.16	(0.27)	0.05
Group III	0.26	(0.24)	0.38	(0.35)	0.13
Tech	0.27	(0.24)	0.41	(0.35)	0.14
Non-Tech	0.19	(0.19)	0.29	(0.29)	0.10
English	0.34	(0.20)	0.52	(0.29)	0.17
Non-English	0.02	(0.08)	0.03	(0.12)	0.01

Positive *change* means level of inconclusiveness increased. Only data from participants who completed both mid-survey and post-survey is included

- 634 – Half of the users (50 %) changed how much they trust paid Lion militia informants, from
- 635 a mean (and standard deviation) of 5.69 (1.66) to 5.38 (1.45).
- 636 – Five (31 %) people changed their opinion about which group is most violent between the
- 637 Mid-survey and Post-survey, and all changed from Reds to Lions. However, users’ confi-
- 638 dence in their answers to this question declined, from a mean (and standard deviation)
- 639 of 7.19 (1.17) to 6.88 (1.59).

640 The quantitative differences recorded are not statistically significant, nonetheless, it is an
 641 important result that many users changed their opinions as a result of interacting with
 642 ArgTrust.

643 With respect to changed levels of confidence, the differences in confidence levels are
 644 shown in Table 12. Averaged across all users who completed both Mid- and Post-surveys,
 645 the confidence level *decreased* after working with ArgTrust. Even considering averages
 646 within a number of groupings (Tech majors versus Non-Tech majors; native English speakers
 647 versus non-native English speakers; Groups II and III), all the averages showed a decline in
 648 confidence. When looking at individuals, the confidence level only increased for 2 out of 16
 649 participants. Both of these individuals were native English speakers; only one was a Tech
 650 major. This result is unexpected. Again, although the quantitative values are not statistically
 651 significant, nonetheless, it is an important result that users’ confidence levels declined.

652 With respect to changed number of Inconclusive answers to conclusive answers (i.e.,
 653 TRUE or FALSE), results are shown in Table 13. For all users who completed both Mid-
 654 and Post-surveys, the percentage of multiple choice questions that were given Inconclusive
 655 answers *increased* between the Mid-survey and the Post-survey. The same holds true for
 656 each sub-group of users (Tech majors versus Non-Tech majors; from English versus non-
 657 English speaking homes). However, we note that the “Non-English” group was significantly

Author Proof

4

5

Table 14 Post-survey questions assessing users' impressions of ArgTrust

(a)	How difficult was the scenario to understand?
(b)	How well would you say you understood the scenario BEFORE using ArgTrust?
(c)	How well would you say you understood the scenario AFTER using ArgTrust?
(d)	How much did the ArgTrust software help to visualize the scenario?
(e)	How mentally demanding was the scenario BEFORE using ArgTrust?
(f)	How mentally demanding was the scenario AFTER using ArgTrust?
(g)	How hard was it to make a decision?
(h)	How much did the ArgTrust software help with your decision making?
(j)	Did you think the ArgTrust system easy to use?
(k)	Overall, how helpful did you find the ArgTrust system?
(m)	How physically demanding was the session?
(n)	How hurried or rushed was the pace of the session?
(p)	How successful were you in accomplishing what you were asked to do?
(q)	How hard did you have to work to accomplish your level of performance?
(r)	How insecure, discouraged, irritated, stressed, and annoyed were you?

658 less likely to provide an **INCONCLUSIVE** answer. While we cannot draw any statistically
 659 significant conclusions from this observation because our sample size is too small, we believe
 660 this trend warrants further investigation in future studies.

661 4.3.7 Post-survey, Part 2

662 Table 14 lists the fifteen questions contained in Part 2 of the Post-survey that aim to assess
 663 users' impressions of ArgTrust; Table 15 contains the responses. These questions are adapted
 664 from the *NASA Task Load Index (TLX)* [26,27]. Figure 7 plots the mean and standard deviation
 665 computed across all 22 participants in the study. The bars labelled (a)–(r) correspond to the
 666 questions listed in Tables 14 and 15. Discussion follows.

667 Table 16 lists the average changes in *understanding* and *mental demand*, respectively, as
 668 perceived by users before and after interacting with ArgTrust. Overall, the mean is positive
 669 in both cases, indicating that both level of understanding and mental demand *increased* as
 670 a result of using ArgTrust. However, the high standard deviations of both values render the
 671 results inconclusive. Figure 8 plots the changes for each individual who participated in the
 672 study. The values are sorted in ascending order, to preserve the anonymity of participants.
 673 These plots illustrate that the understanding increased for most users. Only 2 participants
 674 reported large decreases in understanding as a result of using ArgTrust, and both of these
 675 were in the Non-Technical group. Approximately a third of participants reported increase in
 676 mental demand as a result of using ArgTrust, though this is understandable because everyone
 677 was new to using the software and it is expected that they would have to think harder when
 678 using new software (versus not using it).

679 The second part of the Post-survey also posed the question:

680 Did you know anything about logical argumentation before participating in this
 681 study?

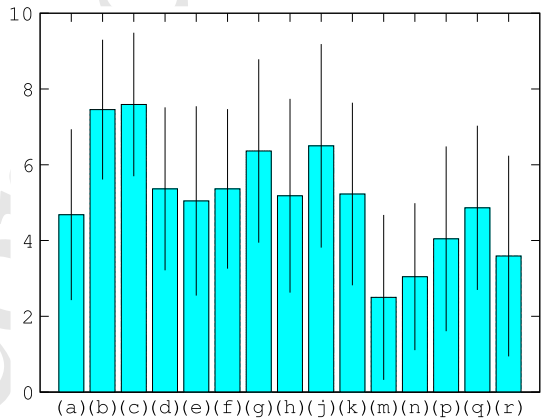
682 which is essentially the same as the following question posed in the Pre-survey:

Table 15 Tallies of post-survey questions assessing users' impressions of ArgTrust

Group	(a)	(b)	(c)	(d)	(e)	(f)	(g)	(h)
Everyone	4.68 (2.25)	7.45 (1.84)	7.59 (1.89)	5.36 (2.15)	5.05 (2.50)	5.36 (2.11)	6.36 (2.42)	5.18 (2.56)
Group I	3.83	7.83	7.00	4.00	5.17	7.00	7.33	4.67
Group II	5.00	7.86	8.57	6.29	5.43	5.00	5.14	6.14
Group III	5.00	6.89	7.22	5.56	4.67	4.56	6.67	4.78
Tech	5.42	7.25	7.75	5.50	4.75	4.83	6.42	4.92
Non-Tech	3.80	7.70	7.40	5.20	5.40	6.00	6.30	5.50
English	4.65	7.47	7.24	5.29	5.18	5.35	6.53	5.06
Non-English	4.80	7.40	8.80	5.60	4.60	5.40	5.80	5.60
Group	(j)	(k)	(m)	(n)	(p)	(q)	(r)	
Everyone	6.50 (2.69)	5.23 (2.41)	2.50 (2.18)	3.05 (1.94)	4.05 (2.44)	4.86 (2.17)	3.59 (2.65)	
Group I	5.50	4.17	4.67	2.67	4.50	5.33	6.33	
Group II	6.00	6.29	1.71	1.86	3.00	5.29	1.71	
Group III	7.56	5.11	1.67	4.22	4.56	4.22	3.22	
Tech	6.58	5.33	1.83	3.75	4.50	4.92	2.92	
Non-Tech	6.40	5.10	3.30	2.20	3.50	4.80	4.40	
English	6.82	4.94	2.18	3.29	4.12	4.71	3.47	
Non-English	5.40	6.20	3.60	2.20	3.80	5.40	4.00	

Responses were given on a scale of 1 (least) to 10 (most)

Fig. 7 Answers to Post-survey assessment questions. The key to the bar labels (a-r) is listed in Table 14

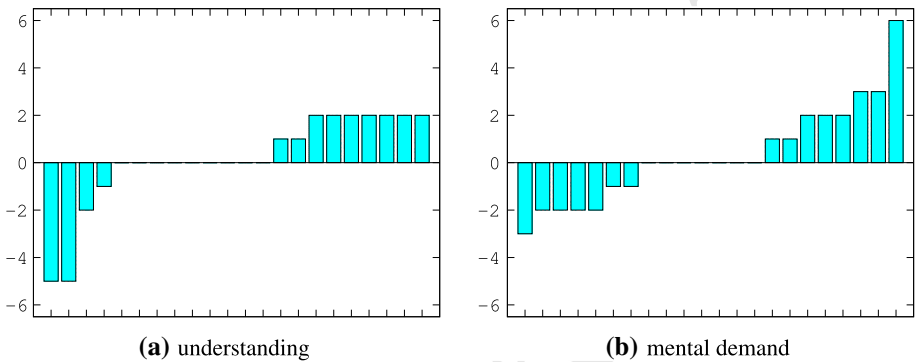


683 Did you know anything about “logical argumentation” or “argumentation graphs”
 684 before today?

685 Table 17 shows the results from both surveys. It is interesting to note that the number of
 686 “yes” responses to questions about knowledge of logical argumentation prior to the study
 687 increased from the Pre-survey to the Post-survey. In fact, 5 people changed their answer from
 688 “no” to “yes” and 2 people changed their answer from “yes” to “no”. Unfortunately, none of

Table 16 Post-survey responses: average self-reported changes in understanding and mental demand

Group	Understanding		Mental demand	
	Mean	SD	Mean	SD
Everyone	0.14	(2.01)	0.32	(2.10)
Group I	-0.83	(3.31)	1.83	(2.71)
Group II	0.71	(0.95)	-0.43	(1.99)
Group III	0.33	(1.41)	-0.11	(1.27)
Tech	0.50	(1.31)	0.08	(1.56)
Non-Tech	-0.30	(2.63)	0.60	(2.67)
English	-0.24	(2.11)	0.18	(2.13)
Non-English	1.40	(0.89)	0.80	(2.17)

**Fig. 8** Post-survey responses: individual self-reported changes in understanding and mental demand. Values are sorted numerically, to preserve anonymity of individuals. Positive values indicate *increases* from “before” to “after”**Table 17** Changes from pre-survey to post-survey with respect to prior knowledge of argumentation

Group	Pre-survey		Post-survey	
	Yes	No	Yes	No
Did you know anything about logical argumentation before participating in this study?				
Everyone	2 (9 %)	20 (91 %)	5 (23 %)	17 (77 %)
Group I	0 (0 %)	6 (100 %)	1 (17 %)	5 (83 %)
Group II	1 (14 %)	6 (86 %)	1 (14 %)	6 (86 %)
Group III	1 (11 %)	8 (89 %)	3 (33 %)	6 (67 %)
Tech	2 (17 %)	10 (83 %)	3 (25 %)	9 (75 %)
Non-Tech	0 (0 %)	10 (100 %)	2 (20 %)	8 (80 %)
English	1 (6 %)	16 (94 %)	5 (29 %)	12 (71 %)
Non-English	1 (20 %)	4 (80 %)	0 (0 %)	5 (100 %)

689 the participants who changed their answer provided any explanation in their comments for
 690 why their answer changed. We speculate that the people who changed their answer from “no”
 691 to “yes” found that they knew about the concept but were unfamiliar with the term “logical
 692 argumentation”, and that the people who changed their answer from “yes” to “no” found that

Table 18 Comments from participants

-
- The tool is successful in representing the factors that go into making a decision and displaying the relationships between those facts and how the influence a decision.
 - It somewhat helped filter out the information that is not necessarily true.
 - It helped me break the component information down a little bit but I had already created an outline of my own that helped me just as much if not more. Seeing the beliefs and the trust broken down was helpful, though.
 - It helped to tell me what I'm supposed to think...how much I'm supposed to trust people, and how I was supposed to interpret the statements in the given scenario. However, it bothered me that a tool was telling me how to simplify a complex problem, since I don't believe the tool can possibly take all the details and subtleties into account. but if I accept that a complex situation can and must be simplified, then yes, the tool is helpful as a place to plug in parameters and let it do the math.
 - By breaking down the situation into smaller bits and displaying how much you believe each situation to be true, it was much easier to make decisions because I was considering the situation asked only, not the entire situation.
 - When I put how I felt into numbers, it organized and simplified my concerns and weighed all of the factors into the equation for me. It made it easier to see the results.
-

693 the concept they thought of as “logical argumentation” was not aligned with the concept as
694 presented in the study.

695 Finally, the Post-survey ends with a free-form comments question, asking users to provide
696 feedback. Table 18 contains some of the free-form comments entered. Though the comments
697 also contain some negative sentiments, we believe that these indicate that some users under-
698 stand the benefits of structuring decisions in the way that we do in ArgTrust.

699 4.4 Analysis

700 In analysing the data, there are a number of interesting observations to make. There was no
701 correlation found between participants' level of education or most other Pre-survey demo-
702 graphics to any of the Mid-survey or Post-survey responses. The factors that did correlate
703 are academic major, split between Tech and Non-Tech majors, and native language, split
704 between native English speakers and non-native English speakers. The results are analysed
705 below.

706 First, we look at the extent to which working with ArgTrust caused users to change their
707 minds.

- 708 – After using ArgTrust, users expressed less trust in paid Lion militia informants.
- 709 – After using ArgTrust, users believed that an attack by the Lion rebel group was less likely.
- 710 – After using ArgTrust, most users believed that the Reds rebel group was the most violent.

711 Each of these findings confirms our working hypothesis that ArgTrust will have impact on
712 users' decision-making processes.

713 Next, we look at differences in users' responses to the fifteen questions in Post-survey
714 Part 2, which aim to measure users' perceptions about the effectiveness of interacting with
715 ArgTrust. While Fig. 7 shows the aggregate data across all users, Figs. 9, 10 and 11 illustrate
716 the results by looking at the users divided into different groupings. While there are no statis-
717 tically significant quantitative differences, there are some interesting qualitative observations
718 that seem important.

Fig. 9 Groups I, II and III. The key to the bar labels (*a-r*) is listed in Table 14

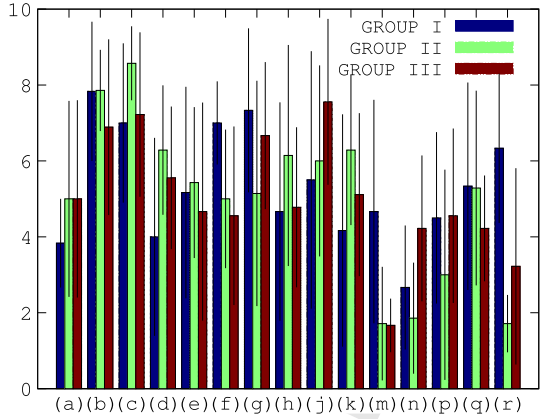


Fig. 10 Technical versus non-technical undergraduate majors. The key to the bar labels (*a-r*) is listed in Table 14

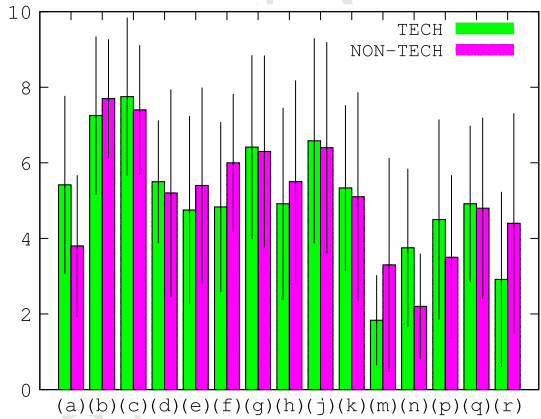
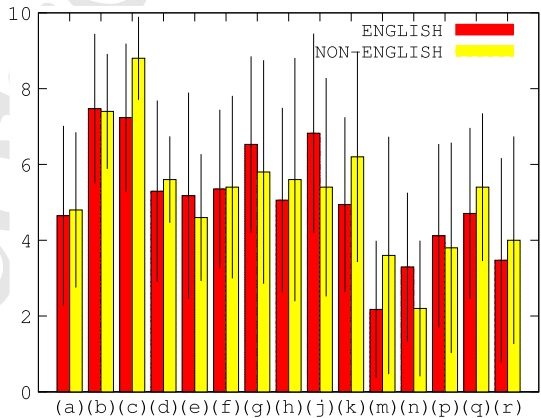


Fig. 11 Participants from English versus non-English speaking homes. The key to the bar labels (*a-r*) is listed in Table 14



719 Specifically, we discuss users' perceived differences between *understanding* and *mental*
 720 *demand* before and after working with ArgTrust. The most distinct differences can be seen
 721 in Fig. 9, bars s and t, respectively.

Author Proof

– With respect to *understanding*, the score for Group II participants rose from 7.86 to 8.57 and the score for Group III participants rose from 6.89 to 7.22. However, Group I reported a decrease in understanding from 7.83 to 7.00. Scores were reported on a scale of 1 (Very poor) to 10 (Very well).

If we look at the numbers of participants reporting an increase/decrease in understanding rather than the average scores, we find that 3 people in Group I reported better understanding and 2 reported worse understanding, 3 people in Group II reported better understanding and none reported worse understanding, and 3 people in Group III reported better understanding and 2 reported worse understanding. Overall, that makes 9 reporting better understanding and 7 reporting worse understanding. Examining the groupings according to Tech majors versus Non-Tech majors, the score went up for Tech majors, but went down slightly for Non-Tech majors.

Examining the groupings according to native English speakers vs non-native English speakers, there is a much greater increase in score for non-native English speakers. difference in users' perceived understanding after working with ArgTrust (c).

– With respect to *mental demand*, the results were more mixed. Groups II and III both showed a decrease in difficulty from pre-software to post-software (5.43–5.00 and 4.67–4.56, respectively). Group I, however, showed a large rise in difficulty (5.17–7.00). Scores were reported on a scale of 1 (Not demanding at all) to 10 (Very demanding).

If we again look at numbers of participants reporting, we find that 1 person from Group I reports the task is less demanding with the software and 4 report it is more demanding, and the figures for Group II and Group III are 4 (less demanding) and 2 (more) and 2 (less) and 2 (more) respectively. Thus, overall, we have 7 reporting it is less demanding and 8 reporting it is more demanding; amongst participants with technical backgrounds, more (6) reported that the task was less demanding when using the ArgTrust software than those who reported the task was more demanding (4).

Next, we examine users' perceived level of frustration when interacting with ArgTrust. Users were asked *How insecure, discouraged, irritated, stressed, and annoyed were you?* in the Post-survey, with responses ranging from 1 (Not annoyed or stressed at all) to 10 (Very annoyed and stressed). The average score for Group II was 1.71, and for Group III the score was 3.22. However, for Group I, the score was significantly higher: 6.33.

Other observations include:

- Group I users did not find the software helpful with respect to visualizing the scenario, as compared to Group II and III users (d).
- Group II found it much easier to make decisions than Group I and III users (g).
- Group III found the ArgTrust system easier to use than Group I and II users (j), while native English speakers found it easier to use than non-native English speakers.
- Non-native English speakers found the system more helpful overall (k).
- For some unknown reason, Group I users found it physically demanding (m) to use the system, which may tie in with their high level of frustration. Otherwise, there was no difference in how the study was implemented; there was no physical element in the study.
- Tech majors felt that they were more rushed or hurried than Non-Tech majors (n).
- Tech majors felt that they were more successful than Non-Tech majors (p).

In summary, the results of the user study provide evidence that participants who are presented with argumentation-based support for making decisions found that interacting with ArgTrust helped their understanding (both directly reported and inferred from the changes made in their decisions), albeit at the cost of increased difficulty in reaching a decision and decreased confidence in their decisions. Nonetheless, users were more able to make

770 decisions after interacting with ArgTrust. Factors that seem to have an impact on users'
771 experience include whether they are native English speakers or not, and whether they majored
772 in Technical subjects as undergraduates or not. Other demographic factors (gender, age,
773 ethnicity, level of education) did not play a part.

774 5 Related work

775 There are four main areas of work that are related to the results reported here: modelling
776 trust; reasoning about trust using argumentation; argumentation-based decision making; and
777 the argumentation in interaction between humans and agents. We briefly cover each of these
778 below stating the differences with our work.

779 5.1 Modelling trust

780 As computer systems have become increasingly distributed, and control of those systems
781 has become more decentralised, computational approaches to trust have become steadily
782 more important [23]. Some of this work has directly been driven by changes in technology,
783 for example considering the trustworthiness of nodes in peer-to-peer networks [1, 15, 32], or
784 dealing with wireless networks [22, 34, 65]. Other research has been driven by changes in the
785 way that technology is used, especially involving the Internet. One early challenge is related
786 to the establishment of trust in e-commerce [47, 62, 79], and the use of reputation systems to
787 enable this trust [36, 37]. Another issue is the problem of deciding which of several competing
788 sources of conflicting information one should trust [2, 11].

789 Additional issues have arisen with the development of the social web, for example, the
790 questions of how social media can be manipulated [41, 42] and how one should revise one's
791 notions of trust based on the past actions of individuals [25]. In this area is some of the work
792 that is most relevant for that we describe here, work that investigates how trust should be
793 propagated through networks of individuals [24, 29, 35, 78], and we have drawn on this in our
794 implementation of ArgTrust.

795 In all of this work, the focus is on the computation of the right numerical trust value to
796 assign to individuals, and is often concerned with showing that the approach advanced in
797 the paper is better than some other approach. In contrast, our work is concerned with how
798 the values computed by these methods can be used in combination with argumentation, and
799 *ArgTrust* provides a choice of ways of computing trust measures drawn from existing models.

800 5.2 Reasoning about trust using argumentation

801 The second area of work to consider is that which looks at the use of argumentation to handle
802 trust. While the literature on trust is considerable, prior work on argumentation and trust
803 is much more sparse. Existing work includes the application of argumentation techniques
804 to networks of trust and distrust [28]; the combination of argumentation with fuzzy trust
805 measures [64], with statistical data [44] and with subjective logic [49] (used in [29] to handle
806 trust); the use of metalevel argumentation to describe trust [73]; and the description of a set
807 of argument schemes for deriving trust [50]. However, none of this covers the same ground
808 as the work reported here which is not concerned with the detail of how argumentation about
809 trust is carried out, but how the results of the argumentation process is used by humans.

810 5.3 Argumentation-based decision making

811 The third area of work to consider is that on argumentation-based decision making. Here
 812 the work by Fox and colleagues [14, 17] showed that constructing arguments for and against
 813 a decision option, and then simply combining these arguments⁸ could provide a decision
 814 mechanism that rivalled the accuracy of probabilistic models. This basic method was extended
 815 in [18, 51] to create a symbolic mechanism that, like classical decision theory, distinguished
 816 between belief in propositions and the values of decision outcomes, while [30] showed the
 817 usefulness of arguments in communicating evidence for decision options to human users.
 818 More recent work on argumentation and decision making is described in [31, 71].

819 The relationship that we consider between argumentation and decision-making is different
 820 from all of this work. All the above work tries to build argumentation systems that identify
 821 the best decision to take. Even [30], which of the work in this area is closest to what we are
 822 doing, tries to identify the best decision and explain it to a human user in terms the human
 823 can understand. In work that evaluates the effectiveness of the systems, the aim is to show
 824 that the system gets the decisions right [14, 77]. Our focus, in contrast, is just to present
 825 information and test whether the users find the information to be useful—whether it helps
 826 them to feel more comfortable with their decisions, and whether they alter their opinion as
 827 a result of being able to use *ArgTrust* to visualise and manipulate the information on which
 828 they base their decisions.

829 5.4 Argumentation for human–agent interaction

830 A number of authors have examined whether software that presents arguments to users can
 831 help their argumentation skills. Work on assessing the impact of software that diagrams argu-
 832 ments has been carried out for several of the better known tools for visualizing arguments.
 833 For example, [10] tests the use of the Questmap system in teaching legal argumentation;
 834 [66] describes formative evaluation of the Belvedere system to support students in scientific
 835 argumentation; [33] examines the use of the TC3 tool to help in collaborative argumentation
 836 in the pursuit of a writing task; and [63] summarises an experiment to assess whether the
 837 “reasoner’s workbench” *Convince Me* helps students, again in the area of scientific reason-
 838 ing. All these papers find some support for the usefulness of the software, but this line of
 839 work also includes [72], which provides a sharp critique of much of the work listed above.
 840 Also noteworthy is [13], which investigates argument mapping—on paper, rather than using
 841 software—as a means for improving understanding of text, and finds that the presentation of
 842 arguments improves comprehension and recall.

843 Another line of related work is that which assesses whether human reasoning matches
 844 that captured in formal argumentation. While work such as [74], which looks at reasoning
 845 with defaults, is somewhat related, the only work we know of that explicitly examines the
 846 differences between argumentation theory and human reasoning is [60]. This paper looks at
 847 the question of reinstatement. Argumentation theory says that if an argument A is attacked
 848 by an argument B and that is all we know then A is OUT and B is IN. However, if B is
 849 then attacked by C , A and C are IN and B is OUT— A is reinstated by the attack that C
 850 makes on B . After reinstatement, the theory says that A is as strongly held as it was before
 851 any attack by B . The study in [60] reveals that while the human subjects understood that A
 852 should be reinstated by the attack from C on B —so they agreed with the general pattern of

⁸ Even the very straightforward mechanism of counting arguments for and against.

853 reinstatement—they did not agree that *A* was held with the same force that it was held before
854 the attack from *B*.

855 Of the work in this final area, [13] is the closest to ours, though differs in its normative
856 focus—it is concerned with whether the students better understand and recall information
857 when it is presented to them using argument diagramming—whereas our work is concerned
858 with how users feel about decisions reached when using argumentation as opposed to deci-
859 sions reached given only text.

860 6 Summary

861 We have described *ArgTrust*, an interactive application designed to help users balance infor-
862 mation from multiple sources and draw conclusions from that information, using logical
863 argumentation and a computational model of trust in information sources. The underlying
864 inference engine is an implementation of earlier work with ArgTrust, here developed using a
865 MySQL and Python back-end and a web-based front-end, largely written in PHP and HTML/CSS.
866 This allows more flexibility with scenario input and the ability to transition to dynamic sce-
867 narios, one of the next steps with this research. As well, this architecture offers more options
868 for user interface development.

869 The main contribution of this paper is the presentation and in-depth analysis of a user
870 study, conducted to evaluate the effectiveness of ArgTrust in helping human users to make
871 decisions. Users were presented with an ambiguous and complex scenario and instructed to
872 formulate an evidence-backed recommendation for an action to take, based on information
873 provided in the scenario. Participants completed Pre-, Mid- and Post-surveys that provided
874 data about their backgrounds (demographics), and their understanding of the scenario before
875 and after interacting with ArgTrust.

876 The results demonstrate, most importantly, that ArgTrust helped users consider their deci-
877 sions more carefully. Analysis of the results provided other qualitative observations: differ-
878 ences in users' experiences with the application, with respect to difficulty understanding the
879 scenario and effort expended in making decisions, were more closely aligned to the lan-
880 guage users speak at home (English or not) and users' undergraduate majors (Technical or
881 not), versus other demographic factors. An unexpected result is that users' questioned their
882 answers more, after working with ArgTrust, indicated by their level of confidence in answers
883 declining after using the application.

884 Future work involves extending the software to handle dynamic scenarios, applying the
885 methodology to additional domains and testing with a wider range of users.

886 **Acknowledgments** This research was funded under Army Research Laboratory Cooperative Agreement
887 Number W911NF-09-2-0053, by the National Science Foundation under grant #1117761, and by the National
888 Security Agency under the Science of Security Lablet grant (SoSL). Additional funding was provided by a
889 University of Liverpool Research Fellowship and by a Fulbright-King's College London Scholar Award. The
890 views and conclusions contained in this document are those of the authors and should not be interpreted as
891 representing the official policies, either expressed or implied, of the funders. The U.S. Government is authorized
892 to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

893 References

- 894 1. Abrams, Z., McGrew, R., & Plotkin, S. (2004). Keeping peers honest in EigenTrust. In *Proceedings of*
895 *the 2nd Workshop on the Economics of Peer-to-Peer Systems*
896 2. Adler, B. T., & de Alfaro, L. (2007). A content-driven reputation system for the Wikipedia. In *Proceedings*
897 *of the 16th International World Wide Web Conference, Banff, Alberta, May.*

- 898 3. Amgoud, L. (1999). Contribution à l'intégration des préférences dans le raisonnement argumentatif. PhD
899 thesis, Université Paul Sabatier, Toulouse.
- 900 4. Amgoud, L., & Cayrol, C. (2002). A reasoning model based on the production of acceptable arguments.
901 *Annals of Mathematics and Artificial Intelligence*, 34(3), 197–215.
- 902 5. Amgoud, L., Maudet, N., & Parsons, S. (2000). Modelling dialogues using argumentation. *Proceedings*
903 *of the Fourth International Conference on Multi-Agent Systems* (pp. 31–38). Boston, MA: IEEE Press.
- 904 6. Azhar, M. Q., Schneider, E., Salvit, J., Wall, H., & Sklar, E. I. (2013). Evaluation of an argumentation-
905 based dialogue system for human-robot collaboration. In *Proceedings of the Workshop on Autonomous*
906 *Robots and Multirobot Systems (ARMS) at Autonomous Agents and MultiAgent Systems (AAMAS), St*
907 *Paul, MN, USA, May*.
- 908 7. Baroni, P., Caminada, M., & Giacomin, M. (2011). An introduction to argumentation semantics. *The*
909 *Knowledge Engineering Review*, 26, 365–410.
- 910 8. Besnard, P., & Hunter, A. (2001). A logic-based theory of deductive arguments. *Artificial Intelligence*,
911 128, 203–235.
- 912 9. Birnbaum, L., Flowers, M., & McGuire, R. (1980). Towards an AI model of argumentation. In *Proceedings*
913 *of the 1st National Conference on Artificial Intelligence* (pp. 313–315).
- 914 10. Carr, C. S. (2003). Using computer supported argument visualization to teach legal argumentation. In P.
915 A. Kirschner, S. J. Buckingham-Shum, & C. S. Carr (Eds.), *Visualizing argumentation: Software tools*
916 *for collaborative and educational sense-making* (pp. 75–96). London: Springer.
- 917 11. Dong, X. L., Berti-Equille, L., & Srivastava, D. (2009). Integrating conflicting data: The role of source
918 dependence. In *Proceedings of the 35th International Conference on Very Large Databases, Lyon, France,*
919 *August*.
- 920 12. Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning,
921 logic programming and n -person games. *Artificial Intelligence*, 77, 321–357.
- 922 13. Dwyer, C. P., Hogan, M. J., & Stewart, I. (2013). An examination of the effects of argument mapping on
923 students' memory and comprehension performance. *Thinking Skills and Creativity*, 8, 11–24.
- 924 14. Emery, J., Walton, R., Coulson, A., Glasspool, D., Ziebland, S., & Fox, J. (1999). Computer support
925 for recording and interpreting family histories of breast and ovarian cancer in primary care (RAGs):
926 Qualitative evaluation with simulated patients. *British Medical Journal*, 319(7201), 32–36.
- 927 15. Feldman, M., Papadimitriou, C., Chuang, J., & Stoica, I. (2004). Free-riding and whitewashing in Peer-
928 Peer systems. In *Proceedings of the 3rd Annual Workshop on Economics and Information Security*.
- 929 16. Ferrando, S. P., & Onaindia, E. (2012). Defeasible argumentation for multi-agent planning in ambient
930 intelligence applications. In V. Conitzer, W. van der Hoek, L. Padgham, & M. Winikoff (Eds.), *Proceed-*
931 *ings of the 11th International Conference on Autonomous Agents and Multiagent Systems*. IFAAMAS:
932 Valencia, Spain.
- 933 17. Fox, J., Glowinski, A., Gordon, C., Hajnal, S., & O'Neil, M. (1990). Logic engineering for knowledge
934 engineering: Design and implementation of the oxford system of medicine. *Artificial Intelligence in*
935 *Medicine*, 2(6), 323–339.
- 936 18. Fox, J., & Parsons, S. (1998). Arguing about beliefs and actions. In A. Hunter & S. Parsons (Eds.),
937 *Applications of uncertainty formalisms*. Berlin: Springer-Verlag.
- 938 19. García, A. J., & Simari, G. (2004). Defeasible logic programming: An argumentative approach. *Theory*
939 *and Practice of Logic Programming*, 4(1), 95–138.
- 940 20. Golbeck, J. (2005). Computing and applying trust in web-based social networks. PhD thesis, University
941 of Maryland, College Park.
- 942 21. Golbeck, J. (May 2006). Combining provenance with trust in social networks for semantic web content
943 filtering. In *Proceedings of the International Provenance and Annotation Workshop, Chicago, Illinois*.
- 944 22. Govindan, K., Mohapatra, P., & Abdelzaher, T. F. (2010, December). Trustworthy wireless networks:
945 Issues and applications. In *Proceedings of the International Symposium on Electronic System Design,*
946 *Bhubaneswar, India*.
- 947 23. Grandison, T., & Sloman, M. (2000). A survey of trust in internet applications. *IEEE Communications*
948 *Surveys and Tutorials*, 4(4), 2–16.
- 949 24. Guha, R., Kumar, R., Raghavan, P., & Tomkins, A. (2004). Propagation of trust and distrust. In *Proceedings*
950 *of the 13th International Conference on the World Wide Web*.
- 951 25. Hang, C.-W., Wang, Y., & Singh, M. P. (2008). An adaptive probabilistic trust model and its evaluation. In
952 *Proceedings of the 7th International Conference on Autonomous Agents and Multiagent Systems, Estoril,*
953 *Portugal*.
- 954 26. Hart, S. G. (2006). NASA-task load index (NASA-TLX); 20 years later. *Proceedings of the Human*
955 *Factors and Ergonomics Society Annual Meeting*, 50, 904–908.
- 956 27. Hart, S. G., & Staveland, L. E. (1988). Development of NASA-TLX (task load index): Results of empirical
957 and theoretical research. *Advances in Psychology*, 52, 139–183.

- 958 28. Harwood, W. T., Clark, J. A., & Jacob, J. L. (2010). Networks of trust and distrust: Towards logical
 959 reputation systems. In D. M. Gabbay & L. van der Torre (Eds.), *Logics in Security, Copenhagen, Denmark*.
 960 29. Jøsang, A., Hayward, R., & Pope, S. (2006). Trust network analysis with subjective logic. In *Proceedings*
 961 *of the 29th Australasian Computer Society Conference, Hobart, January*.
 962 30. Judson, P. N., Fox, J., & Krause, P. J. (1996). Using new reasoning technology in chemical information
 963 systems. *Journal of Chemical Information and Computer Sciences*, 36, 621–624.
 964 31. Kakas, A., & Moraitis, P. (2003). Argumentation based decision making for autonomous agents. In *2nd*
 965 *International Conference on Autonomous Agents and Multi-Agent Systems*. New York, NY: ACM Press.
 966 32. Kamvar, S. D., Schlosser, M. T., & Garcia-Molina, H. (2004). The EigenTrust algorithm for reputation
 967 management in P2P networks. In *Proceedings of the 12th World Wide Web Conference, May*.
 968 33. Kanselaar, G., Erkens, G., Andriessen, J., Prangma, M., Veerman, A., & Jaspers, J. (2003). Designing
 969 argumentation tools for collaborative learning. In P. A. Kirschner, S. J. Buckingham-Shum, & C. S. Carr
 970 (Eds.), *Visualizing argumentation: Software tools for collaborative and educational sense-making* (pp.
 971 51–73). London: Springer.
 972 34. Karlof, C., & Wagner, D. (2003). Secure routing in wireless sensor networks: Attacks and countermea-
 973 sures. *Ad Hoc Network*, 1, 293–315.
 974 35. Katz, Y., & Golbeck, J. (2006). Social network-based trust in prioritized default logic. In *Proceedings of*
 975 *the 21st National Conference on Artificial Intelligence*.
 976 36. Khopkar, T., Li, X., & Resnick, P. (2005). Self-selection, slipping, salvaging, slacking and stoning: The
 977 impacts of. In *Proceedings of the 6th ACM Conference on Electronic Commerce, June*. Vancouver: ACM.
 978 37. Khosravifar, B., Bentahar, J., Moazin, A., & Thiran, P. (2010). On the reputation of agent-based web
 979 services. *Proceedings of the 24th AAAI Conference on Artificial Intelligence, July* (pp. 1352–1357).
 980 Atlanta: AAAI Press.
 981 38. Kirschner, P. A., Buckingham-Shum, S. J., & Carr, C. S. (Eds.). (2003). *Using computer supported*
 982 *argument visualization to teach legal argumentation*. Berlin: Springer.
 983 39. Kok, E., Meyer, J.-J., Prakken, H., & Vreeswijk, G. (2012). Testing the benefits of structured argumentation
 984 in multi-agent deliberation dialogues. In *Proceedings of the 9th International Workshop on Argumentation*
 985 *in Multiagent Systems, Valencia, Spain*.
 986 40. Kraus, S., Sycara, K., & Evenchik, A. (1998). Reaching agreements through argumentation: A logical
 987 model and implementation. *Artificial Intelligence*, 104(1–2), 1–69.
 988 41. Lang, J., Spear, M., & Wu, S. F. (2010). Social manipulation of online recommender systems. In *Pro-*
 989 *ceedings of the 2nd International Conference on Social Informatics, Laxenburg, Austria*.
 990 42. Lerman, K., & Galstyan, A. (2008). Analysis of social voting patterns on Digg. In *Proceedings of the 1st*
 991 *Workshop on Online Social Networks, Seattle, August*.
 992 43. Liau, C.-J. (2003). Belief, information acquisition, and trust in multi-agent systems—a modal logic for-
 993 mulation. *Artificial Intelligence*, 149, 31–60.
 994 44. Matt, P.-A., Morge, M., & Toni, F. (2010). Combining statistics and arguments to compute trust. In
 995 *Proceedings of the 9th International Conference on Autonomous Agents and Multiagents Systems, Toronto,*
 996 *Canada, May*.
 997 45. Mercier, H., & Sperber, D. (2011). Why do humans reason? Arguments for an argumentative theory.
 998 *Behavioral and Brain Sciences*, 34(2), 57–74.
 999 46. Modgil, S., & Prakken, H. (2013). A general account of argumentation with preferences. *Artificial Intel-*
 1000 *ligence*, 195, 361–397.
 1001 47. Mui, L., Moteashemi, M., & Halberstadt, A. (2002). A computational model of trust and reputation. In
 1002 *Proceedings of the 35th Hawai'i International Conference on System Sciences*.
 1003 48. Naylor, S. (2005). *Not a good day day to die: The untold story of operation Anaconda*. New York: Berkley
 1004 Caliber Books.
 1005 49. Oren, N., Norman, T., & Preece, A. (2007). Subjective logic and arguing with evidence. *Artificial Intel-*
 1006 *ligence*, 171(10–15), 838–854.
 1007 50. Parsons, S., Atkinson, K., Li, Z., McBurney, P., Sklar, E., Singh, M., et al. (2014). Argument schemes for
 1008 reasoning about trust. *Argument and Computation*, 5(2–3), 160–190.
 1009 51. Parsons, S., & Green, S. (1999). Argumentation and qualitative decision making. In *Proceedings of the*
 1010 *5th European Conference on Symbolic and Quantitative Approaches to Reasoning with Uncertainty*.
 1011 52. Parsons, S., McBurney, P., & Sklar, E. (May 2010). Reasoning about trust using argumentation: A position
 1012 paper. In *Proceedings of the Workshop on Argumentation in Multiagent Systems, Toronto, Canada*.
 1013 53. Parsons, S., Sierra, C., & Jennings, N. R. (1998). Agents that reason and negotiate by arguing. *Journal of*
 1014 *Logic and Computation*, 8(3), 261–292.
 1015 54. Parsons, S., Sklar, E. I., Salvit, J., Wall, H., & Li, Z. (2013). ArgTrust: Decision making with information
 1016 from sources of varying trustworthiness (Demonstration). In *Proceedings of Autonomous Agents and*
 1017 *Multiagent Systems (AAMAS), St Paul, MN, USA, May*.

- 1018 55. Parsons, S., Tang, Y., Sklar, E., McBurney, P., & Cai, K. (2011). Argumentation-based reasoning in agents
1019 with varying degrees of trust. In *Proceedings of the 10th International Conference on Autonomous Agents*
1020 *and Multi-Agent Systems, Taipei, Taiwan*.
- 1021 56. Parsons, S., Wooldridge, M., & Amgoud, L. (2003). Properties and complexity of formal inter-agent
1022 dialogues. *Journal of Logic and Computation*, 13(3), 347–376.
- 1023 57. Pasquier, P., Hollands, R., Rahwan, I., Dignum, F., & Sonenberg, L. (2011). An empirical study of
1024 interest-based negotiation. *Journal of Autonomous Agents and Multi-Agent Systems*, 22(2), 249–288.
- 1025 58. Prakken, H. (2000). On dialogue systems with speech acts, arguments, and counterarguments. In *Pro-*
1026 *ceedings of the Seventh European Workshop on Logic in Artificial Intelligence*. Berlin: Springer Verlag.
- 1027 59. Prakken, H. (2005). Coherence and flexibility in dialogue games for argumentation. *Journal of Logic and*
1028 *Computation*, 15, 1009–1040.
- 1029 60. Rahwan, I., Madakkatel, M. I., Bonnefon, J. F., Awan, R. N., & Abdallah, S. (2010). Behavioral exper-
1030 iments for assessing the abstract argumentation semantics of reinstatement. *Cognitive Science*, 34(8),
1031 1483–1502.
- 1032 61. Rahwan, I., & Simari, G. R. (Eds.). (2009). *Argumentation in Artificial Intelligence*. Berlin: Springer
1033 Verlag.
- 1034 62. Resnick, P., & Zeckhauser, R. (2002). Trust among strangers in internet transactions: Empirical analysis
1035 of eBay's reputation system. In M. R. Baye (Ed.), *The economics of the internet and E-commerce* (pp.
1036 127–157). Amsterdam: Elsevier Science.
- 1037 63. Schank, P., & Ranney, M. (1995). Improved reasoning with Convince Me. In *CHI'95 Conference Com-*
1038 *panion* (pp. 276–277).
- 1039 64. Stranders, R., de Weerd, M., & Witteveen, C. (2008). Fuzzy argumentation for trust. In F. Sadri & K.
1040 Satoh, (Eds.), *Proceedings of the Eighth Workshop on Computational Logic in Multi-Agent Systems*.
1041 Lecture Notes in Computer Science (vol. 5056, pp. 214–230). Berlin: Springer Verlag.
- 1042 65. Sun, Y., Yu, W., Han, Z., & Liu, K. J. R. (2005). Trust modeling and evaluation in ad hoc networks. In
1043 *Proceedings of the Yth Annual IEEE Global Communications Conference* (pp. 1862–1867).
- 1044 66. Suthers, D., Weiner, A., Connelly, J., & Paolucci, M. (1995). Belvedere: Engaging students in critical
1045 discussion of science and public policy issues. In *Proceedings of the 7th World Conference on Artificial*
1046 *Intelligence in Education, Washington, DC, August* (pp. 266–273).
- 1047 67. Sycara, K. (1990). Persuasive argumentation in negotiation. *Theory and Decision*, 28, 203–242.
- 1048 68. Tang, Y., Cai, K., McBurney, P., Sklar, E., & Parsons, S. (2012). Using argumentation to reason about
1049 trust and belief. *Journal of Logic and Computation*, 22(5), 979–1018.
- 1050 69. Tang, Y., Cai, K., Sklar, E., & Parsons, S. (2011). A prototype system for argumentation-based rea-
1051 soning about trust. In *Proceedings of the 9th European Workshop on Multiagent Systems, Maastricht,*
1052 *Netherlands, November*.
- 1053 70. Tang, Y., Sklar, E. I., & Parsons, S. (2012). An argumentation engine: ArgTrust. In *Proceedings of the*
1054 *Workshop on Argumentation in Multiagent Systems (ArgMAS) at Autonomous Agents and MultiAgent*
1055 *Systems (AAMAS), Valencia, Spain, June*.
- 1056 71. Tolchinsky, P., Modgil, S., Cortes, U., & Sanchez-Marre, M. (2006). Cbr and argument schemes for
1057 collaborative decision making. In *Proceedings of the First International Conference on Computational*
1058 *Models of Argument, Liverpool* (pp. 71–82).
- 1059 72. van den Braak, S. W., van Oostendorp, H., Prakken, H., & Vreeswijk, G. A. (2006). A critical review
1060 of argument visualization tools: Do users become better reasoners? In *Proceedings of the Workshop on*
1061 *Computational Models of Natural Argument* (pp. 67–75).
- 1062 73. Villata, S., Boella, G., Gabbay, D. M., & van der Torre, L. (2011). Arguing about the trustworthiness
1063 of the information sources. In *Proceedings of the European Conference on Symbolic and Quantitative*
1064 *Approaches to Reasoning and Uncertainty, Belfast, UK*.
- 1065 74. Vogel, C. M. (1995). Inheritance reasoning: Psychological plausibility, proof theory and semantics. PhD
1066 thesis, University of Edinburgh, Centre for Cognitive Science.
- 1067 75. Vreeswijk, G., & Prakken, H. (2000). Credulous and sceptical argument games for preferred semantics.
1068 In *Proceedings of the 7th European Workshop on Logics in Artificial Intelligence*.
- 1069 76. Walton, D. N., & Krabbe, E. C. W. (1995). *Commitment in dialogue: Basic concepts of interpersonal*
1070 *reasoning*. Albany, NY: State University of New York Press.
- 1071 77. Walton, R., Gierl, C., Mistry, H., Vessey, M. P., & Fox, J. (1997). Evaluation of computer support for
1072 prescribing (CAPSULE) using simulated cases. *British Medical Journal*, 315, 791–795.
- 1073 78. Wang, Y., & Singh, M. P. (2006). Trust representation and aggregation in a distributed agent system. In
1074 *Proceedings of the 21st National Conference on Artificial Intelligence, Boston, MA*.
- 1075 79. Yu, B., & Singh, M. (2002). Distributed reputation management for electronic commerce. *Computational*
1076 *Intelligence*, 18(4), 349–535.