

Reasoning with Knights and Knaves:  
Towards an understanding of reasoning about  
truth and falsity

Floris van Vugt  
floris.van.vugt@ens.fr

January 31, 2009

1 **1 Introduction**

2 At the end of the eighties a type of logical puzzle called *knight and knave* prob-  
3 lems(Smullyan, 1987) enters the scene of psychology of reasoning(Rips, 1989).  
4 They are staged on an imaginary island where only two kinds of people live:  
5 knights, who always tell the truth, and knaves, who always lie. It is furthermore  
6 assumed that they have complete and correct knowledge of each other's being  
7 knight or knave. A puzzle then consists of a number of utterances from a few of  
8 these inhabitants. The task for the reader is to decide of each character whether  
9 he or she is a knight or a knave.

10 To illustrate, consider the example in table 1.

11 As with all puzzles, there are multiple ways to arrive at the solution. One

Table 1: A sample problem from (Johnson-Laird & Byrne, 1990)

- A: A and B are knaves.
- B: A is a knave.

12 could start out making the assumption that A is a knight. This means both  
13 A and B must be knaves, but that is contrary to our assumption. If however  
14 one would have assumed A a knave, then he must be lying, i.e. not both A and  
15 B are knaves. Since he himself is a knave by assumption, that leaves as only  
16 possibility that B is knight. And indeed what B says is true. Since one of our  
17 two assumptions led to a contradiction, we conclude the other is correct and  
18 that A is a knave and B knight.

## 19 **1.1 Problem solvability**

20 First of all, as none of the existing literature has made explicit, the above puzzle  
21 is rather unique in that it has one and exactly one solution. As a matter of fact  
22 the *solvable* puzzles can be said to reside on a thin line in between what one  
23 could call *paradoxical* problems on the one hand and *underspecified* problems  
24 on the other.

### 25 **1.1.1 Paradoxical problems**

26 I will call a problem *paradoxical* if any attribution of knight- and knave-status  
27 to the speakers leads to a contradiction. For example,

- 28 • A: B is a knight.
- 29 • B: A is a knave.

30 Clearly if A is a knight, then B must be a knight, but at the same time B  
31 must then be lying in saying that A is a knave, so there is a contradiction. If A  
32 is a knave then he must lie and therefore B is a knave, however what B says is  
33 suddenly true.

34 An even more primitive example of such a paradoxical phrase is the direct  
35 translation of the Liar sentence (Kripke, 1975)(Tarski, 1983) which says of itself  
36 that it is false,

- 37 • A: A is a knave.

### 38 1.1.2 Underspecified problems

39 A problem is *underspecified* if there are multiple attributions of knight- and  
40 knave-status to the speakers that are consistent. For example,

41 • A: B is a knight.

42 • B: A is a knight.

43 If A is a knight, then she must be telling the truth, hence B is a knight also,  
44 which is consistent with what he says. If A is a knave however, then she must  
45 be lying and B therefore is a knave, indeed it is a lie that A is a knight.

46 The problem is that this sequence cannot be solved on the basis of the  
47 utterances we have.

48 Again, the most primitive form of such sentences is found in Truth-teller  
49 sentences(Kripke, 1975)(Tarski, 1983) and its equivalent on the island of knights  
50 and knaves would be

51 • A: A is a knight

## 52 1.2 Outline

53 The purpose of this paper is to provide a very modest overview of the discussion  
54 that took place from the early nineties onwards between the main players in the  
55 psychology of reasoning and which revolved around the knight-and-knave prob-  
56 lems. Also, my aim will be to provide my personal reflections on the arguments  
57 presented.

## 58 2 Rips and the mental rules approach

59 Rips(Rips, 1989) was the first to suggest these brain teasers as an object of study  
60 for the psychology of reasoning. His motivation is that so far the field has focused  
61 on a very narrow body of reasoning tasks such as Aristotelian syllogisms and  
62 Wason's selection task. However, one could assert that his more or less hidden  
63 agenda was to address the question to what extent psychological theories should

Table 2: Rips' knight–knave–specific rules

1. if  $\text{says}(x, p)$ ,  $\text{knight}(x)$  then  $p$
2. if  $\text{says}(x, p)$ ,  $\text{knave}(x)$  then  $\neg p$
3. if  $\neg\text{knave}(x)$  then  $\text{knight}(x)$
4. if  $\neg\text{knight}(x)$  then  $\text{knave}(x)$

64 appeal to semantic concepts such as truth and falsity (Rips, 1986). Rips' hope  
65 is to demonstrate that his inference rule–based model sufficiently explains how  
66 a subject handles these problems without explicitly requiring a notion of truth  
67 or falsity on the level of the theory. This would in a broader context serve as  
68 an argument that in cognitive science such semantic concepts are superfluous.

## 69 2.1 The approach of mental deduction rules

70 On the basis of an informal observation of subjects solving these puzzles, Rips  
71 suggests the following model for their reasoning.

72 The model considers mental deduction rules a psychological primitive and  
73 they are used to calculate conclusions from a limited number of assumptions.  
74 The core propositional rules are adapted from a generic model (Rips, 1983) (Rips,  
75 1989) that is capable of performing elementary inferences. This general proposi-  
76 tional model is then supplemented with knight–knave–specific deduction rules,  
77 which represent the content of the instructions that one gives to the subject,  
78 e.g. that knights always speak the truth and knaves always lie. These specific  
79 rules are listed in table 2.<sup>1</sup>

80 Given these rules, the strategy for solving the puzzles can be represented as  
81 the following computer program.

- 82 1. It begins by assuming that the first speaker is a knight.

---

<sup>1</sup>In what follows,  $\text{says}(x, p)$  will represent that  $x$  utters “ $p$ ,” and  $\neg p$  is the negation of  $p$ ,  
i.e. NOT  $p$ .

83 2. From this assumption and using the generic and specific deduction rules,  
84 the program derives as many conclusions as possible. This phase stops  
85 when either of the following obtain:

86 (a) The set of assumptions and conclusions is inconsistent. In this case  
87 the assumption that the first speaker is a knight is abandoned and  
88 replaced by the assumption that he is a knave.

89 (b) No more rules apply, i.e. none of the deduction rules yields a conclu-  
90 sion that was not already found. In this case the program proceeds  
91 with the assumption that the second speaker is a knight.

92 3. The program continues like this until it has found all consistent sets of  
93 assumptions about the knight- and knave-status of the individuals.

94 The essential statement of the model is that the total number of applications  
95 of the rules needed to arrive at a conclusion is a measure for the complexity of  
96 the problem. This means that it predicts that problems that require a larger  
97 number of steps will take subjects longer and they will make more errors on  
98 them. In order to eliminate the influence of irrelevant factors, Rips forms pairs  
99 of problems that contain the same number of speakers and clauses, i.e. atomic  
100 propositions, but require different numbers of steps to solve them. By comparing  
101 subject performance within each pair only, Rips thus cancels out influence of  
102 processes other than actually solving the puzzle, such as reading the problem  
103 statement.

#### 104 **2.1.1 Suppositional reasoning**

105 It is interesting to note that the program is *suppositional*, that is, it starts out  
106 by making an assumption, an *Ansatz*. Rips decided this was the most authentic  
107 procedure since he observed in an informal experiment where subjects were  
108 instructed to think aloud while solving the problems that they all started by  
109 making such an assumption and seeing where that reasoning led. However, when  
110 this procedure was reproduced in follow-up studies reference to the amount

Table 3: A model structure in the rule-based derivation.

A	B
knight	?

111 of suppositions is at least ambiguous. For instance (Elqayam, 2003) explains  
 112 subjects make a supposition about the status of the first speaker and derive  
 113 its consequences, and then make the contrary assumption, and “[t]hey thus  
 114 proceed” (p.268). This means, they continue to make the supposition the second  
 115 speaker is a knight, and then that he is a knave, and then the same for the  
 116 third and so on. This is the only correct way to interpret Rips’ explanation of  
 117 the procedure (Rips, 1989) (p.91). Moreover it is not difficult to see that many  
 118 problems would not even be solvable without making such multiple assumptions.

119 What is remarkable is that in each consistent set (Rips, 1989) (p.91), *ever*  
 120 speaker thus is once the object of an assumption, not just the first speaker.  
 121 This, in my point of view calls the question to what extent the subject in Rips’  
 122 model is not actually constructing *models* and using the natural deduction rules  
 123 to verify that they are consistent. Rips’ derivation requires the subject to have  
 124 at every point in time least some sort of a structure which can be visualised as  
 125 a table and which keeps track of what speakers have what status.

126 For instance, in the problem mentioned in the introduction, the subject,  
 127 after the assumption that A is a knight he needs a structure like in table 3 to  
 128 represent that he has made one assumption about A and none yet about B.

129 That is to say, I argue that if Rips’ mental rules are a psychological primitive,  
 130 at least some sort of mental model is also.

### 131 2.1.2 Origin of the knight–knave rules

132 The model proceeds by repeated *rule application*. A first question could be why  
 133 subjects would have precisely these rules in mind. Particularly puzzling is the  
 134 absence of *backward inference rules* like

135 if says( $x, p$ ) and  $p$  then knight( $x$ )

136 The reason they are not in the model is that Rips' bases his model on the  
137 subjects' thinking aloud when solving the problems and he never observed them  
138 using such *backward inference rules*(Rips, 1989)(p.89).

139 Moreover Rips' model does not need these rules, since the procedure of  
140 reasoning outlined above will simulate its behaviour. For instance, if says( $x, p$ )  
141 and  $p$ , then the program will begin by assuming  $\text{knight}(x)$ , which is precisely the  
142 conclusion we wanted. If it would consider the opposite, i.e.  $\text{knave}(x)$ , then it  
143 runs into a contradiction upon application of the knave-rule 2 of table 2 because  
144 it yields  $\neg p$ .

145 However, personally I found myself using this rule directly in a number of  
146 problems. The response that Rips could give here is that the number of inference  
147 steps his model yields is nevertheless a *measure* of how long it takes a subject  
148 to solve the problem even though the subject might occasionally optimise his  
149 strategy since such heuristics can probably be applied across the problem types  
150 equally.

151 Rather disconcerting is that Rips later in the article finds himself forced to  
152 add additional rules to his model to allow it to solve a wider variety of knight-  
153 knave-problems. Indeed one might long for a proof that the given model is  
154 capable of solving all problems, unless of course one wants to allow for the  
155 possibility that there are formally solvable problems that no human can solve.

### 156 **2.1.3 Determinism of rule application**

157 Another interesting note that the literature has not picked up is that Rips'  
158 procedure relies on the decidability of the deduction rules. This is crucial in the  
159 second step of the program, where as many conclusions as possible are drawn  
160 from the previously discovered or assumed facts. If this procedure would not  
161 eventually terminate in a state where the application of any rule no longer  
162 leads to a conclusion that was not already drawn, the subject would continue  
163 to derive new conclusions without ever passing to different assumptions. That  
164 is, the program would not be guaranteed to halt.

165 It seems plausible that the rules Rips proposes have this decidability prop-  
166 erty, even though he does not explicitly prove it. The reason is that he formu-  
167 lates uniquely elimination rules which only reduce the already finite complexity  
168 of phrases in the set of conclusions and therefore the procedure will eventually  
169 run out(Gentzen, 1969).

170 What remains remarkable, however, is that subjects would have precisely  
171 such a set of decidable rules in their minds. There are numerous examples  
172 in logic where different axiomatisations which are equivalent in terms of con-  
173 clusions that are derivable from them, nevertheless differ in their decidability.  
174 One of these examples is Lambek’s application of Gentzen’s sequent calculus to  
175 phrase structures(Lambek, 1958) where he starts with the most intuitive set of  
176 rules, which is not decidable, and then needs several important revisions before  
177 arriving at a decidable set with the same proof-theoretic power.

178 If Rips is right that reasoning works with mental deduction rules, we are  
179 then also faced with the question why, of all possible axiomatisations, we have  
180 a decidable one in our heads.

#### 181 **2.1.4 Optimisation in problem solving**

182 Rips devotes a single sentence to report that his program also uses an opti-  
183 malisation heuristic: “After each step, the program revises the ordering of its  
184 rules so that rules that have successfully applied will be tried first of the next  
185 round.”(Rips, 1989)(p.91). But this raises the question how Rips’ model ac-  
186 counts for this structure in subjects’ performance if he wants to do away with  
187 any metalogical reasoning. For to be able to decide on the “success” of a rule it  
188 seems one needs a certain representation of one’s one reasoning in the previous  
189 step.

190 However, if again we grant Rips that such optimisations can be performed  
191 equally well across the experimental conditions their effect (i.e. the lowering of  
192 the number of steps required for solution) will be overall and therefore cancel  
193 out when comparing subject’s performance in different problems.

## 194 **2.2 Experimental confirmation of mental rules–account**

195 In a first experiment Rips registers only the accuracy of the subject’s responses.  
196 First of all he observed widespread incapacity to solve the problems: 10 out of 34  
197 subjects gave up on the experiment within 15 minutes, and among the subjects  
198 that completed the test solved on average only about 20% of the problems was  
199 correctly answered. Second, among the pairs of problems matched for number  
200 of clauses and speakers more errors are made on the more difficult ones.

201 In a second experiment he measured the time it takes subjects to solve  
202 two–speaker three–clause problems. Again the main finding is that in spite of  
203 high error rates, subjects take longer to solve problems that take model more  
204 inference steps to solve.

## 205 **3 Critique of mental rules and introduction of** 206 **mental models**

### 207 **3.1 Evans**

208 The first critique of Rips’ study comes from (Evans, 1990).

209 First of all, Evans argues, the knight and knave problems are not meaningful  
210 in the real–world context, where one hardly encounters people who either always  
211 lie or speak the truth. This means in particular that it is doubtful to what extent  
212 subjects’ performance in the experiment reflects reasoning as it is employed *in*  
213 *vivo*: “[w]e must recognise that almost all real–world cognition occurs in the  
214 presence of meaningful context” (Evans, 1990)(p.86–87).

215 Secondly, Evans feels the procedure Rips proposes for solving these riddles  
216 is unjustifiably deterministic in the sense that it eventually always finds the  
217 correct answer. The observations with towering high error rates contain only a  
218 fraction of correct responses, so in the best possible scenario Rips’ model can  
219 be applied to this fraction only. The errors themselves can hardly be accounted  
220 for by a model that has no way to “generate” these errors itself.

221 Finally, though he himself places less emphasis on it, he nevertheless raises  
222 the interesting remark that Rips' model takes as a starting point the puzzle  
223 encoded in a logical format, e.g.  $\text{says}(x, p \wedge q)$  rather than "X says that p and  
224 q." Although Evans, nor any other author that I know of, for that matter,  
225 does not develop this further, it does point into a seemingly trivial but essential  
226 nuance that might not be clear from the problem description: the scope of  
227 conjunction. For instance, the natural language version of our example could  
228 also have been transcribed as  $\text{says}(x, p) \wedge \text{says}(x, q)$ . This distinction is crucial  
229 for it turns the problem into a completely different one and I would even go as  
230 far as to argue that at least part of the errors can be attributed to this kind  
231 of misunderstanding. For instance, the problem in table 1 becomes paradoxical  
232 as soon as we would take A as uttering two assertions which therefore need to  
233 both be true or both be false. The first of which,  $\neg A$  would render the puzzle  
234 paradoxical.

## 235 **3.2 Johnson–Laird and Byrne**

236 The next substantial criticism comes from a hardly surprising corner (Johnson-  
237 Laird & Byrne, 1990).

### 238 **3.2.1 Criticism of the mental deduction rules**

239 The main problem Johnson–Laird and Byrne identify in Rips' approach is again  
240 the deterministic nature of the procedure he describes. On the one hand it seems  
241 unrealistic to assume that subjects come to the task with a ready-made solution  
242 procedure as effective as Rips' model, and on the other hand it seems that even  
243 if they had, the procedure is so powerful that it would place unrealistic demands  
244 on their computational facilities.

245 Essentially, this problem lies in the need to follow up on disjunctive sets of  
246 models. For instance, if one speaker asserts  $p \wedge q$  and the program arrived at the  
247 point of assuming that speaker a knight, it will then have to follow-up on each of  
248 the situations  $\{p, \neg q\}$ ,  $\{\neg p, q\}$  and  $\{\neg p, \neg q\}$  and especially when it would need

249 to compute additional disjunctive situations concerning other speakers in each  
250 case, the number of cases to be considered would grow exponentially, placing  
251 impossible demands on subject's memory.

252 However, as Rips argues in his defense(Rips, 1990), Johnson-Laird and  
253 Byrne seem to have misrepresented his position although they claim to use  
254 a simple "notational variant"(p.73). Though Rips does not explain this further,  
255 most probably he refers to the fact that his program will never consider such  
256 disjunctive cases separately but simply derive whatever conclusion is possible  
257 from the statement of the disjunction. In the example of  $\neg(p \wedge q)$  the identities  
258 in his natural deduction model lead to conclude  $\neg p \vee \neg q$  and then leave it at  
259 that.

260 A point Rips himself did not raise but which seems equally valid, is that even  
261 if subjects were in some way required to compute these disjunctive cases, then  
262 perhaps these "impossible demands" are precisely an explanation of the high  
263 error rates. To develop this further, one would need to determine which of the  
264 presented problems required following-up on disjunctive sets and see whether  
265 they yielded higher error rates and reaction times.

### 266 **3.2.2 Model approach: developing strategies**

267 Their point of view is that reasoning is based on *mental models*, or "internal  
268 model[s] of the state of affairs that the premises describe."(Johnson-Laird &  
269 Byrne, 1991)(p.35). Instead of deriving conclusions using rules without neces-  
270 sarily knowing what sorts of situations, or extensions, these conclusions refer  
271 to, Johnson-Laird and Byrne propose that reasoning is the construction and  
272 manipulation of mental representations that are more or less explicit. Broadly  
273 speaking, when a subject performs a modus ponens, he or she starts with a men-  
274 tal representation in which both premises are verified and then tries to create  
275 a model in which these remain true but the conclusion is false. Once he or she  
276 realises this cannot be done, the modus ponens is accepted as logically valid.

277 They feel it unreasonable to assume that subjects already have a ready-made

278 procedure for solving the puzzles and rather develop ways, called *strategies*, to  
279 solve them as they observe themselves working.

280 They suggest to account for the data observed by Rips as the workings of  
281 four such mental strategies that are much like heuristics and which result from  
282 subjects' observing themselves performing the task: "With experience of the  
283 puzzles, they are likely to develop more systematic strategies." (Johnson-Laird  
284 & Byrne, 1990)(p.72). This is the kind of meta-cognitive capacity they feel Rips  
285 tried to evade in his model.

286 The proposed strategies are the following:

287 1. *Simple chain*. This strategy is to, like in Rips' model, follow all the conse-  
288 quences of assuming the first speaker to be a knight, with one difference:  
289 once one is required to look into disjunctive consequences, that is, pre-  
290 cisely the case described before, which they identified as problematic in  
291 terms of cognitive complexity.

292 2. *Circular*. Once a speaker utters something that is self-referential, such  
293 as: "I am a knave and B is a knave," then the strategy is to follow up  
294 only on the immediate consequences, i.e. those that require a single rea-  
295 soning step, since those often already rule out one of the cases. Thus  
296 the strategy is to not pursue the consequences of the consequences. In  
297 our example, assuming that the speaker is a knight can in such a way be  
298 rejected instantly.

299 3. *Hypothesise-and-match*. This strategy involves matching other speaker's  
300 utterances to previous conclusions. For instance, consider the following  
301 example:

- 302 • A: A and B are knights.
- 303 • B: A is a knave.

304 The point is that as soon as one concludes that A cannot be a knight, and  
305 therefore must be a knave, then we can match this conclusion with B's  
306 assertion. Since they are the same thing, B must be a knight.

307 Interestingly, this deduction is precisely the inverse rule mentioned in sec-  
308 tion 2.1.2.

309 4. *Same-assertion-and-match*. In the case where two speakers make the  
310 same assertion, any other speaker who attributes a different status to  
311 them is necessarily lying.

- 312 • A: C is a knave.
- 313 • B: C is a knave.
- 314 • C: A is a knight and B is a knave.

315 In a post-hoc analysis of Rips' very own data, they then proceed to show  
316 that problems which can be solved using these four strategies yield significantly  
317 more correct answers than those who cannot.

### 318 **3.2.3 Reflection on mental models account**

319 I would like to remark that the *simple chain* and *circular* strategies (and possibly  
320 the other two as well) only serve to eliminate parts of the "tree" of cases to be  
321 considered for a complete solution. As such, they are what in information science  
322 would be called a *heuristic*, they cut down parts of the search tree but they do  
323 not alter significantly the nature of the problem solution.

324 Secondly, there appears to be no unsystematic theory that unites them and  
325 therefore they can be said to be *ad hoc* in the sense that it would not be a  
326 surprise if one would come up with another strategy or maybe conclude that  
327 one of them is not applied after all. The problem about this is that the model  
328 has too many free parameters and therefore escapes scientific testing, rendering  
329 it pseudoscientific in a Popperian sense. Equivalently, it is very doubtful what  
330 the strategies really *explain* in subject's performance.

### 331 **3.2.4 Rips' response to mental strategies**

332 Rips responds in considerable detail(Rips, 1990) to the criticism outlined before.

333       Hardly surprisingly, one of his first remarks is that there is not much that  
334 the *mental models* contribute to Johnson–Laird and Byrne’s approach to the  
335 problem. The strategies could have been formulated equally easy in a mental  
336 deduction rule framework, as in one based on mental models. Therefore, first  
337 of all, they do not particularly confirm the mental models account as such.

338       Furthermore, Rips remarks that in their post–hoc analysis of his data, of the  
339 four strategies, the *circular* was not included in the test for any puzzle in which  
340 it applied could have also been solved by the simple chain. Similarly, the *same–*  
341 *assertion–and–match* strategy because it applied in too little cases to allow  
342 statistical comparison. Then, if one matches the problems for number of clauses  
343 and speakers of the remaining two only *hypothesise–and–match* significantly  
344 explain the difference in scores. In other words, there is only experimental  
345 evidence for one of the four strategies.

346       However, in a study focussing more broadly on strategies in reasoning, Byrne  
347 and Handley(Byrne & Handley, 1997), mounting experiments of their own, find  
348 further evidence for reasoning strategies, taking away much of the power of this  
349 objection of Rips’.

350       Finally, in a remarkably lucid passage that, unfortunately to my knowledge  
351 has not been followed up in the literature, Rips also clarifies his position con-  
352 cerning the rejection of the use of meta–logical notions in psychological theories.  
353 Although Johnson–Laird and Byrne and Evans for that matter have taken him  
354 to reject using the notion of truth altogether, he argues only against appealing  
355 to expert theories of truth to explain subject’s behaviour. Rips feels a theory  
356 can call on stage the subject’s representation of truth, but it should not go fur-  
357 ther than that by using some independent theory of truth that logicians provide  
358 us with to explain how subjects behave: “Although cognitive psychologists can  
359 investigate people’s beliefs about truth... it is quite another thing for cognitive  
360 psychologists to explain behaviour by appeal to the nature of truth itself.” (Rips,  
361 1990)(p.296–297)

Table 4: Example utterance from (Elqayam, 2003)(p.280)

- I am a knave or I am a knight

## 362 4 Elqayam and the norm in knight and knave 363 puzzles

364 At this point in time the discussion between Rips, Evans and Johnson–Laird and  
365 Byrne falls quiet. More than a decade later, new light is shed on the discussion  
366 by Shira Elqayam(Elqayam, 2003).

### 367 4.1 Truth–value gaps

368 Among her most profound comments is that so far all studies into the knight and  
369 knave puzzles have assumed that there is a single “correct” answer. However,  
370 Elqayam observes the knight–knave puzzles presented to the subjects contained  
371 instances of the Liar and Truth–teller sentences.

372 These sentences are the starting point of Kripke’s theory of truth, because  
373 they show that one cannot define a truth–predicate such that “ $p$  is true” is  
374 true if and only if “ $p$ .”(Kripke, 1975) From there onwards several solutions are  
375 proposed, most of them introducing a third, “undefined” truth value in addition  
376 to “true” or “false,” or, equivalently, a true predicate simply not applying to a  
377 certain number of sentences, like the liar.

378 Consider for instance the utterance in table 4.

379 Since the island is supposed to contain only knights and knaves, one can  
380 consider this phrase a tautology. On the other hand, Elqayam argues it can  
381 equally well be considered false since neither of the subphrases is necessarily true  
382 and some authors in philosophical logic classify such phrases as false. Finally,  
383 as long as the knight– or knave–status of the speaker has not been determined  
384 we can consider the two subphrases as undefined, i.e. the third truth value, and  
385 hence also their disjunction. Thus, depending on the norm we apply one can

386 justifiedly consider a phrase either true, false, or neither.

387 This directly undermines the definition of a “correct” answer and thus might  
388 provide an essential clue as to the nature of the large number of “errors” ob-  
389 served. She argues that this absence of an objective norm could be remedied  
390 by allowing subjects when they classify speakers as either knight or knave the  
391 option that they “do not know.”

## 392 **4.2 Reflection on truth–gaps**

393 I think Elqayam’s observation of the implicit assumption of a logical norm in  
394 computing the response “correctness” is invaluable, and deserves as much credit  
395 as David Hume whom Immanuel Kant thanked for rousing him from his dog-  
396 matic slumber.

397 Before introducing my criticism, I would like to point out that Elqayam is  
398 precisely embarking in the analyses that Rips warned against (Rips, 1990) that  
399 is, she appeals to expert theories of truth to explain subject’s behaviour. I would  
400 agree with Rips in the sense that it is important not to take the truth theories  
401 as restricting the possibilities of reasoning. Instead, Elqayam’s analysis appears  
402 to me valid in that it hypothesises what notion of truth the subject uses when  
403 solving the task.

### 404 **4.2.1 Truth–value gaps violate an instruction**

405 The instruction given to the subjects stating that each inhabitant of the island  
406 is either knight or knave, is equivalent to the law of excluded middle. Therefore,  
407 the explicit instruction to the subject is to operate in bivalent logic. That is, as  
408 soon as a subject would consider that what a certain speaker has said is neither  
409 true nor false, he has violated the aspect of the puzzle that every inhabitant is  
410 either knight or knave and therefore in a way he or she is no longer solving the  
411 puzzle that was originally given. Thus, although Elqayam might offer a valid  
412 explanation of the “errors” observed in Rips’ original experiment, it is not an  
413 example of the subject “justifiedly” using a different norm, which is what she

414 argues.

415 At this point it is interesting to notice firstly the parallel with children solving  
416 the Tower of Hanoi problem, where often they are observed impose themselves  
417 additional constraints<sup>2</sup>. The difference here is that if subjects consider truth-  
418 value gaps in knight-knave puzzles they not elaborated the puzzle but they  
419 simply ignored one of its essential instructions: the law of excluded middle.

#### 420 4.2.2 Paradoxality

421 In response to Elqayam’s observation, it is good to remind ourselves that none of  
422 the speakers refers *only* to himself in their utterances. In those cases the problem  
423 could have also been formulated by eliminating that utterance, because at best  
424 they are redundant by not adding anything to the problem and at worst they  
425 cause the problem to be “underspecified,” to use the distinction I introduced  
426 before. For instance, the puzzle of table 4 was never part of a problem presented  
427 to the subjects. This means that in particular, sentences such as the Liar and  
428 the Truth-teller, which are so far the only compelling reasons for us to abandon  
429 a bivalent truth assignment, do not occur.

430 It seems that Elqayam has confounded self-referentiality with paradoxality.  
431 This has been recently a greatly investigated topic in logic. Broadly speaking,  
432 Yablo showed an example of a paradox without self-reference (Yablo, 1993) and,  
433 conversely, Leitgeb argues in a recent paper that many sentences that refer  
434 to themselves can be considered not paradoxical (Leitgeb, 2005). The example  
435 Elqayam gives herself also falls in this latter category. In conjunction, these  
436 results show that paradoxality and self-referentiality are far from being the  
437 same thing. Ironically, Elqayam seems to have applied a high-level version of the  
438 *circularity* strategy of Johnson-Laird and Byrne, suspecting paradox as soon as  
439 a speaker refers to himself.

---

<sup>2</sup>They take many more steps to solve the Tower of Hanoi problem since they do not allow themselves to move a stone two piles away.

440 **4.2.3 Paradox by circularity and paradox by excluded–middle**

441 Let us then turn to sentences which *could* and *did* occur in the problems pre-  
442 sented to the subjects and look a bit closer at why they would contain a truth-  
443 value gap. For instance,

- 444 • A: I am a knight or B is a knave.

445 Elqayam would consider the first part (“I am a knight”) as undetermined, in  
446 analogy to the liar sentence, which is undetermined. The reason is most likely  
447 that she feels there is a certain circularity analogous to the liar sentence, where,  
448 if we want to know whether it is true or false, we first need to know whether it  
449 itself is true or false, thus begging the question.

450 In modern logic and especially in recent days there has been considerable  
451 research into this idea, called *groundedness*(Leitgeb, 2005). The idea is that to  
452 determine the truth or falsity of certain sentences, like “It is true that snow  
453 is white,” one needs to know the truth or falsity of “Snow is white” and that  
454 sentence itself does not depend on another sentence but on a state of affairs in  
455 the external world of which we are capable of verifying whether it is the case.  
456 Therefore, knowing this state of affairs we can fill in the truth value of “It is true  
457 that snow is white.” This is why we tend to consider such sentences *grounded*.

458 However, in order to know the truth or falsity of a sentence like “This sen-  
459 tence is false.” we would need to first know whether the sentence itself is true,  
460 for which we need to look at sentence itself again, and so on infinitely. This  
461 vicious circularity is why we call such sentences *ungrounded*.

462 And precisely here dawns a very important distinction between knight–knave  
463 puzzles and truth–predicate definition: in the latter case liar sentences are para-  
464 doxical because of *circularity* (for sentences become true or false by virtue of  
465 what they express being the case or not), in the former because of the knight–  
466 knave–island variant of the *excluded middle*.

467 If an inhabitant of the knight–knave island utters: “I am a knave,” then  
468 in that will force us to abandon the assumption that all inhabitants are either  
469 knight or knave, if at all we want to evade contradiction. In that respect, even

470 switching to trivalent logic would not help. But if an inhabitant utters: “What  
471 I now say is false,” *that* will force us to abandon bivalent logic and with it also  
472 conclude that the one who utters it is neither knight nor knave.

473 Put in another way, we assume that each inhabitant is either a knight or a  
474 knave, even before he or she has said anything. The inhabitant does not *become*  
475 knight or knave by the uttering of a truth or a lie, he or she is assumed to have  
476 been so all along. It is only to *us*, listeners and explorers of the island, that  
477 their status turns from “indeterminate” *for us* to knight or to knave.

478 Thus, when Elqayam praises Rips for including a “do not know” option in  
479 his first experiment or other researchers(Schroyens *et al.* , 1999) for including  
480 even response patterns reminiscent of four-valued logic(Gupta & Belnap, 1993),  
481 that does not point subjects to three- or four-valued logic, but simply expresses  
482 their incapacity to tell.

#### 483 **4.2.4 Three-valued-logic and suppositional reasoning**

484 The merit of Elqayam’s proposal of the application of multivalued logic in the  
485 knight-knave puzzles has thus brought to light an essential difference between  
486 the knight-knave puzzles and the definition of a truth predicate in logic. The  
487 difference is that the island of knights and knaves, it seems, contains an ad-  
488 ditional layer where truth-value gaps can appear. For instance, if a person is  
489 neither knight or knave that would make for a “local” truth-value gap that  
490 violates the instruction that each person is either knight or knave. If a person  
491 utters a liar sentence, however, that makes for a “global” truth-value gap that  
492 violates bivalent logic.

493 Also, the distinction between paradox by circularity and paradox by excluded-  
494 middle helps to understand why the solution procedures proposed by all authors  
495 dealing with the knight-knave puzzles so far have always been *suppositional* (see  
496 section 2.1.1). That is, Rips already observed subjects need to start out by sup-  
497 posing a speaker to be either knight or knave and then deduce consequences.  
498 The point is that only making the supposition a speaker is a knight and then

499 the supposition that the speaker is a knave will reveal the paradox by excluded-  
500 middle, whereas a paradox by circularity will yield a contradiction already by  
501 application of deduction rules. For instance, the liar sentence is shown to be  
502 paradoxical as soon as one substitutes it in the Tarski T-equivalence “ $p$  is true”  
503 iff  $p$ .

#### 504 4.2.5 Bivalent logic

505 So where do we go then, if switching to trivalent logic does not help to explain  
506 the outcome of Rips’ original experiment?

507 A clue might come from one of the most influential papers in contemporary  
508 logic (Leitgeb, 2005). Leitgeb proposes a definition of a predicate of truth which  
509 evades paradoxes while remaining in two-valued logic. This is achieved by  
510 applying the “naive” condition for truth predicates<sup>3</sup> only to grounded sentences.  
511 The unique feature of this approach to logical paradox that stays within bivalent  
512 logic and seems therefore the most appropriate candidate to handle knight-  
513 knave puzzles where the excluded-middle principle is an explicit constraint.

514 It would be interesting to use knight and knave puzzles to test whether sub-  
515 jects actually use such a conception of truth. Like Rips (Rips, 1990) emphasised,  
516 “[t]here is also no doubt that people have common-sense beliefs about truth and  
517 falsity, and it is of interest to document these notions and to compare them with  
518 expert theories.” Perhaps, using the knight-knave paradigm, this question can  
519 actually be brought into the realm of experimental verification.

520 My very modest proposal is to eliminate the instruction that all inhabitants  
521 are either knight or knave. Thus, the only thing we instruct the subjects is that  
522 knights always tell the truth and knaves always lie.

523 Then consider the problems in table 5. The idea is that even though A utters  
524 an ungrounded sentence, B could be said to be a knight in virtue of knowing  
525 that snow is white or that a person cannot both be a knight and a knave.

---

<sup>3</sup>That is, the Tarski T-equivalence that a sentence “ $p$  is true” is true if and only if “ $p$ ” is true

Table 5: Testing a subject’s conception of truth

Problem I

- A: I am a knave.
- B: A is a knight or snow is white.
- Puzzle: What is B?

Problem II

- A: I am a knave.
- B: A is not both a knight and a knave.
- Puzzle: What is B?

526 If subjects turn out to be able to solve these two problems, one can conclude  
527 that the law of excluded middle is not inherent in their reasoning. For if it were,  
528 they would run aground upon hearing what A says. If subjects are not able  
529 to solve these problems that would corroborate Elqayam’s point that subjects  
530 reason using a trivalent logic.

531 I realise there are many problems with this task and it is quite beyond  
532 the scope of this paper to deal with them. My aim was mainly to point out  
533 the possibility that knight–knave puzzles can help to understand how subjects  
534 conceive truth, and perhaps in the future inspire a more thoughtful analysis.

## 535 **5 Conclusion**

536 We have seen almost two decades of research into how subjects reason to solve  
537 knight–knave brain–teasers. Rips proposed a model based on mental deduction  
538 rules in which we, as psychologists of reasoning, do not need to appeal to meta–  
539 cognition. The results were criticised by Evans and Johnson–Laird and Byrne  
540 who propose their own interpretation based on mental models and meta–logical

541 reasoning strategies.

542 Elqayam, almost a decade later, calls into doubt the nature of the norm  
543 that the previous authors have presupposed to be the only meaningful norm  
544 in knight–knave puzzles. In particular, she argues the problems call for or at  
545 least justify the use of three–valued logic. My commentary is that knight–knave  
546 puzzles come with the explicit requirement of the excluded middle, which forced  
547 us to conclude that subjects who use three–valued logic are no longer solving  
548 the puzzle as it was proposed. This is perhaps the most truthful explanation of  
549 Rips’ observation of high error rates.

550 On the other hand, perhaps more importantly, these considerations can lead  
551 us to view these puzzles in a different way: rather as a tool that might lead to  
552 discover what subjects’ conceptions about truth and falsity are.

## References

- Byrne, Ruth M. J., & Handley, Simon J. 1997. Reasoning strategies for suppositional deductions. *Cognition*, **62**(1), 1 – 49.
- Elqayam, Shira. 2003. Norm, error, and the structure of rationality: The case study of the knight–knave paradigm. *Semiotica*, **2003**(147), 265–289.
- Evans, Jonathan St. B. T. 1990. Reasoning with knights and knaves: A discussion of rips. *Cognition*, **36**(1), 85 – 90.
- Gentzen, Gerhard. 1969. *The collected papers of gerhard gentzen*. M. e. szabo edn. Studies in logic and the foundations of mathematics. Amsterdam: North-Holland Pub. Co.
- Gupta, Anil, & Belnap, Nuel. 1993. *The revision theory of truth*. Cambridge, Massachusetts, and London, England: MIT Press.
- Johnson-Laird, P. N., & Byrne, R. M. 1990. Meta-logical problems: knights, knaves, and rips. *Cognition*, **36**(1), 69–84; discussion 85–90.

- Johnson-Laird, P.N., & Byrne, R.M.J. 1991. *Deduction*. Essays in cognitive psychology. Hove (UK) Hillsdale, NJ (USA): L. Erlbaum Associates.
- Kripke, Saul. 1975. Outline of a theory of truth. *The journal of philosophy*, **72**(19), 690–716.
- Lambek, J. 1958. The mathematics of sentence structure. *American mathematical monthly*, **65**, 154–170.
- Leitgeb, Hannes. 2005. What truth depends on. *Journal of philosophical logic*, **34**, 155–192.
- Rips, L. J. 1989. The psychology of knights and knaves. *Cognition*, **31**(2), 85–116.
- Rips, L. J. 1990. Paralogical reasoning: Evans, Johnson-Laird, and Byrne on liar and truth-teller puzzles. *Cognition*, **36**(3), 291–314.
- Rips, Lance J. 1983. Cognitive processes in propositional reasoning. *Psychological review*, **90**(Jan), 38–71.
- Rips, Lance J. 1986. *The representation of knowledge and belief*. Tucson: University of Arizona Press. Myles Brand and Robert M. Harnish. Chap. Mental Muddles.
- Schroyens, W., Schaeken, W., & D’Ydewalle, G. 1999. Error and bias in metapropositional reasoning: A case of the mental model theory. *Thinking and reasoning*, **5**(38), 29–66.
- Smullyan, R.M. 1987. *What is the name of this book? the riddle of Dracula and other logical puzzles*. NJ: Prentice-Hall: Englewood Cliffs.
- Tarski, Alfred. 1983. The concept of truth in formalized languages. In: Corcoran, J. (ed), *Logic, semantics and metamathematics*. Indianapolis: Hackett Publishing Company. The English translation of Tarski’s 1936 *Der Wahrheitsbegriff in den formalisierten Sprachen*.
- Yablo, Stephen. 1993. Paradox without self-reference. *Analysis*, **53**(4), 251–252.