# What to Believe When Inferences are Contradicted: The Impact of Knowledge Type and Inference Rule

**Renée Elio** `(ree@cs.ualberta.ca)`
Department of Computing Science
University of Alberta
Edmonton, Alberta T6G 2H1

## Abstract

Simple belief-revision tasks were defined by a giving subjects a conditional premise, ($p$—>$q$), a categorical premise, ($p$, for a modus-ponens belief-set, or ~$q$, for a modus tollens belief-set), and the associated inference ($q$ or ~$p$, respectively). "New" information contradicted the initial inference (~$q$ or $p$, respectively). Subjects indicated their degree of belief in the conditional premise and the categorical premise, given the contradiction. Results indicated that the choice was a function of the knowledge type expressed in the conditional form; when that knowledge type was causal, the choice was affected by the number of disabling factors associated with the causal relationship. A "possible worlds" interpretation of the data is related to formal notions such as epistemic entrenchment, used in normative models of belief revision, and to reasoning from uncertain premises, from the human deduction literature.

## Introduction

Suppose you initially consider the following to be true: (a) If Jeremy mows the Robinsons' lawn, they will pay him $15; (b) He mowed their lawn; and so, by inference, (c) The Robinsons paid him $15. You then discover that they did <u>not</u> pay him $15. Which of your initial beliefs do you chose to abandon? That he mowed the lawn? Or that they'd give him $15 if he did so?

This example illustrates a problem associated with belief revision, the process by which a reasoner makes the transition from one belief state to another. Sometimes new information will conflict with beliefs currently held to be true. Understanding the principles by which such a conflict is resolved is relevant to certain formal efforts in AI as well as to the psychological community that studies human reasoning. From the AI perspective, there is much work on formalizing "rationality postulates" for belief revision (see Alchourrón, Gärdenfors, & Makinson, 1985). Yet even that community has acknowledged that "extra-logical pragmatic principles are needed to guide the revision process" (Nebel, 1989). Much of the formal AI work on belief revision makes an appeal to an idea called *epistemic entrenchment* (Gärdenfors, 1988) The rationale behind epistemic entrenchment is that, practically, an agent may need to choose among alternative ways to change its beliefs, and intuitively put, some beliefs will be more "deserving" of

continued belief than other beliefs, in the face of contradiction. While epistemic entrenchment is offered as a formal concept, there seems to be little understanding about what the principles for epistemic entrenchment actually are.

From the psychological perspective, human performance on classical deductive problems has been extensively studied. But there is a need for descriptive data and theories on how people resolve inconsistency when new information (about a static world) is obtained. The studies presented in this article are part of a larger research effort concerned exactly with this issue: what principles are involved in how people decide to resolve contradiction? Put another way, which of several former beliefs is called into question when a valid inference that follows from them is subsequently contradicted?

Some AI belief-revision theorists have argued that conditional statements like $p$—>$q$ may warrant, *a priori,* a higher degree of entrenchment, not because there is something to be preferred about material implications, but because that form often signals "law-like" or predictive relations that have explanatory power (e.g., Foo & Rao, 1988). Elio and Pelletier (in press; 1994) did not find any evidence that people were more apt to retain conditional over non-conditional premises, when considering abstract belief-revision problems. In fact, they found that *disbelieving* a conditional statement was the preferred belief-revision decision when the same problems were instantiated with natural-language concepts. Elio and Pelletier applied a kind of theory formation account to their finding that people preferred to disbelieve a conditional ($p$—>$q$) rather than categorical ($p$) statement, as a way of reconciling the new contradictory information. Under this account, a conditional can be viewed as expressing a regularity about the world. If an inference that follows from such a regularity is subsequently contradicted by some new data, then (an inductive perspective might argue) the regularity itself must be flawed. The idea that data enjoys a priority over regularities has been offered as a belief revision principle in other frameworks (e.g., Thagard, 1988) particularly when regularities are (merely) hypotheses under consideration to explain or systematize observed facts.

Elio and Pelletier (in press) identified, but did not address, a broader question of epistemic entrenchment, namely whether different *types* of regularities (expressed as a conditional) might be differentially "entrenched" in the face of contradictory evidence. This indeed seems plausible, for

it has been long recognized (see Evans, Newstead, & Bryne, 1993) that the material implication view of conditionals has proven inadequate from both a linguistic and psychological viewpoint. Natural-language uses of the conditional include definitions, promises, advice, threats, and causal relationships. Within the realm of causal reasoning, the work directly relevant to the present study's design is that done by Cummins and her colleagues (Cummins, Lubart, Alksnis, & Wrist, 1991; Cummins, 1995). They report that deductive reasoning about causality is affected by the number of alternative causes of the consequent and the number of disabling conditions — factors that prevent effects from occurring even in the presence of viable causes. The possible import of this latter distinction for belief revision, as defined here, seems clear: belief revision is required exactly when that which is expected to be true (or false) is not. More generally, the belief-revision "problem"—defined here as reconciling contradictions to deductively valid inferences—has a direct relevance to theoretical accounts of reasoning from uncertain premises. I consider this connection more fully in the Summary section.

The studies reported here examined whether belief-revision decisions are influenced by (a) the type of knowledge expressed in an initially-believed conditional and (b) the type of inference rule (modus ponens—MP—or modus tollens—MT) used to define the initial belief set. The rationale for the former issue has already been outlined. Concerning the MP v. MT question, Elio & Pelletier (1994; in press) have reported that MP and MT belief sets are revised differently. It is important to extend those results to cases in which the kind of knowledge expressed in the conditional is systematically manipulated.

## Inference Rule and Knowledge Types
### The Conditionals
A set of 65 conditionals was constructed to include (a) the 16 causal conditionals used by Cummins *et al*. (1991) and (b) other conditional statements that expressed familiar definitions, unfamiliar definitions, and promises. To ensure that the experimenter's intuitions matched those of subjects, a group of 40 subjects provided classification ratings for each statement, presented as follows:

"If a substance is nitroglycerin, then its molecular structure is $C_3H_5(NO_3)_3$."
I would classify what this sentence is expressing as ....
    a. a promise
    b. a cause and effect
    c. a prediction (not based on cause and effect)
    d. a definition that I am familiar with
    e. a definition that I am not familiar with

(The choices included the category of "prediction" in case subjects believed there was a temporal, but non-causal, contingency between a conditional's antecedent and consequent.) From the ratings received, 34 statements were selected to be used in the belief revision problems. Examples of each type are: (a) *Causal-Many Alternative Causes/Many Disablers*: If the brake was depressed, then the car slowed down; (b) *Causal-Many Alternative Causes/Few Disablers:* If read without his glasses, then he got a headache; (c) *Causal-Few Alternative Causes/Many Disablers*: If the trigger was pulled, then the gun fired; (d) *Causal-Few Alternative Causes/Few Disablers*: If Joe cut his finger, then it bled. (e) *Promise:* If Jeremy mows the lawn, then the Robinsons will give him $15; (f) *Unfamiliar Definition*: If an organism is, then it is a parasite of vertebrate animals; (g) *Familiar Definition:* If a flower is an annual, then it dies after one year of blooming. The causal conditionals were taken from Cummins (1995) and the full set used for this study is given in Elio (1996). Due to space constraints, details on the data collected for definition belief-revision problems will not be presented here.

### The Task
Table 1 shows an example belief revision problem. For the modus ponens problems, the initial belief set consisted of a conditional statement $p—>q$, the categorical statement $p$, and a statement asserting $q$ as a consequence that follows from the other two statements. The update sentence, $\sim q$, contradicted the inference in the initial belief set. For the modus tollens problem, the initial belief set was $p—>q$, $\sim q$, therefore $\sim p$. The update sentence was $p$. The terms "data statement" and "categorical statement" are used interchangeably to refer to the initial belief $p$ in the MP case and $\sim q$ in the MT case.

---

This is what you initially believe:
    If water was poured on the campfire, then the campfire went out.
    Water was poured on the campfire.
    From this, you believe the campfire went out.

You do further investigation and discover:
        The campfire did <u>not</u> go out.

Assuming the new information is true, what do you think the degree of belief should be for.....
       Water was poured on the campfire.
Disbelieve                  Believe
   A———B———C———D———E———F———G
          Uncertain
and for....
 If water was poured on the campfire, then the campfire went out.
Disbelieve                  Believe
   A———B———C———D———E———F———G
          Uncertain

Table 1
Example MP Belief Revision Problem

---

Two types of tasks were used with this material. In both tasks, subjects were instructed to assume that the initial sentences were true. The results reported here concern a degree-of-belief rating task, for which subjects indicated a new degree-of-belief in the initial-belief sentences, *given* the contradiction introduced by the update sentence. The other type of task used was a forced-choice task, following Elio & Pelletier (1994; in press), in which other subjects explicitly indicated which of the initial beliefs they would retain and

which they would abandon. Where relevant, I indicate how the forced-choice results support interpretations based on degree-of-belief results; see Elio (1996) for more detail on the forced-choice data.

**Design, Subjects, and Procedure**
The complete problem set of 34 belief-revision was considered too long for subjects to work through, without risk of boredom. Thus, subjects were randomly assigned to receive one of two booklets. In one booklet type, each of 34 subjects saw all 16 causal conditionals and all six promise conditionals. In the other booklet type, each of 34 subjects saw all 12 definitions and the six promise conditionals. This allowed causal type and definition-type to serve as a repeated factors, for these two groups, respectively. Analyses of the promise problems used data from both these two groups combined, to increase the sample size. The order of the problems and the order in which initial beliefs appeared for rating was randomized for each subject.

**Collecting Prior Belief Ratings**
Subsequent to the study described above, 22 subjects rated all 34 conditionals on a 1-to-7 scale, where 7 was labeled "Likely to be true" and 1 was labeled "Likely to be false."

**Results**
For the scale given in Table 1, the letters A through G were mapped to numerical values 1 (disbelieve) through 7 (believe). The top portion of Table 2 gives, in parentheses, the average belief-ratings to the conditionals rated in isolation. The main body of the table presents the degree-of-belief ratings for conditional and categorical beliefs, after the contradictory information was presented. As I discuss further in the summary, additional prior belief ratings are needed in addition to the ratings for the conditionals in isolation. However, these averages provide some baseline that demonstrates that subjects' level of belief was affected by the contradiction in the belief-revision task.

For purposes of the present study, the key metric is the where the conditional and the categorical premise are placed on the belief scale relative to each other, and the magnitude of the difference between the ratings assigned to them. These can be taken to indicate the degree to which one of those beliefs is singled out over the other as "suspect", given the contradiction.

**Causal Conditionals.** A 2 (inference rule:MP v. MT) X 2 (alternative causes) X 2(alternative disablers) X 2 (item rated: conditional statement, data statement) analysis of variance, with repeated measures on all three factors revealed a main effect for inference rule, and interactions between alternative-disablers and item-rated (F (1,33) = 14.63, p = .001), alternative-causes and item-rated (F(1,33) = 10.30, p = .003), and inference-rule and item-rated (F(1,33) = 7.42, p = .01). Although the patterns of responses appear different for MP and MT problems as a function of causal-knowledge type, the three-way interaction did not approach significance.

The key effect is the role of disabling factors. The focus is on the direction of the differences between the degree-of-belief in the conditional v. the categorical belief, given the contradiction. When the contradicted belief-set involved a few-disabler conditional, subjects retained a higher degree of belief the conditional than in the categorical belief; when the revision problem involved a many-disabler conditional, the conditional was assigned somewhat lower rating than the categorical belief. The clearest contrast is between the many-disablers/many causes case (a -1.0 difference) v. the few-disablers/few-causes case (a 1.8 difference) for the MP problems.

The absolute values for the degree-of-belief ratings does in many cases hover around the "uncertain" mark; however, my interpretation that the disabler-factor influences *which* initial belief is called into question is supported by results I found with the forced-choice task, in which subjects chose among: (a) the believe the conditional and disbelieve the categorical statement; (b) believe the categorical statement and disbelieve the conditional; (c) be uncertain about both. For many-disabler conditionals, 43% of the subjects to retain the conditional and 33% chose to disbelieve it. But for few-disabler conditionals, 54% of subjects retained the conditional and only 26% chose to disbelieve it.

The effect of syntactic inference rule (MP v. MT belief sets) is consistent with results from previous studies on these types of problems: when considering the contradiction, subjects gave higher belief ratings to conditionals than to data statements on the modus-ponens belief sets, and higher ratings to data statements than to conditionals in the modus-tollens belief sets.

| Modus Ponens Belief Set | | | |
|---|---|---|---|
| Upon learning $\sim q$, the new degree of belief in ... | | | |
| | $p \longrightarrow q$ | $p$ | Difference |
| Causals | | | |
| Many Disablers of p—>q | | | |
| (5.7)   Few Causes of q | 4.1 | 4.3 | -0.2 |
| (5.5)   Many Causes of q | 3.7 | 4.7 | -1.0 |
| Few Disablers of p—>q | | | |
| (6.4)   Few Causes of q | 5.1 | 3.3 | 1.8 |
| (5.6)   Many Causes of q | 4.7 | 3.9 | 0.7 |
| Promises | | | |
| (4.3) | 3.3 | 4.8 | -1.6 |

| Modus Tollens Belief Set | | | |
|---|---|---|---|
| Upon learning $p$, the new degree of belief in.... | | | |
| | $p \longrightarrow q$ | $\sim q$ | Difference |
| Causals | | | |
| Many Disablers of p—>q | | | |
| Few Causes of q | 4.0 | 5.0 | -1.0 |
| Many Causes of q | 3.8 | 5.2 | -1.4 |
| Few Disablers of p—>q | | | |
| Few Causes of q | 4.4 | 4.0 | .4 |
| Many Causes of q | 4.0 | 4.7 | -.7 |
| Promises | 3.3 | 5.4 | -2.1 |

Table 2: (Prior Belief Ratings) and New Belief Ratings —Familiar Domains

**Promises.** The difference between the degree of belief in the promise-conditionals v. the promise-data statement was greater for modus-tollens belief sets (-2.1) than for modus-ponens belief sets (-1.6) ($F(1,67) = 4.85$, p =.031). What is notable, however, is that for both MP and MT belief sets, the direction of the effect is large and in the same direction. The degree-of-belief ratings are consistent with the data on the forced-choice task: only 25% of subjects chose to believe the promise conditional and disbelieve the non-conditional; 52% of them opted to disbelieve the promise conditional and retain belief in the non-conditional. Conditionals expressing promises appear more corrigible— i.e., "easier to disbelieve"—in the face of contradiction than conditionals expressing other types of knowledge.

## Discussion

Clearly, the conditional form itself is too crude to serve even as a heuristic upon which to hang principles of epistemic entrenchment. And so too is knowledge-type distinction itself, given the influence of the few-disablers distinction on how subjects revised belief sets involving causal conditionals. These results offer some insight into what the "extra-logical" preferences guiding belief revision might be, and suggest that they can be formalized in a somewhat domain independent fashion. They suggest a view of epistemic entrenchment modeled as the result of assessing the likelihood of the alternative possible worlds that correspond to different ways of accounting for a contradiction. When there are few-disablers for $p—> q$, then there are few(er) possible worlds in which a disabler could be in effect. It is less likely, then, that current world (to be modeled by a belief state) is one of those worlds in which a disabler is also true. In this case, it may be a more plausible belief-revision decision to retain belief in the conditional statement and question the validity of the categorical statement (what I am call the "data statement"). Put another way, the heuristic that "data has priority" may not always lead to a plausible belief revision decision, given a reasoner's background knowledge.

The few-alternative-causes/few-disablers case can be interpreted as a much tighter, almost definitional relation between $p$ and $q$: there are few other explanations for $p$ when $q$ is true, and few other explanations for why $q$ would fail to hold when $p$ is true. The pattern of data for few-causes/few disabler causals in this study looked remarkably similar to the pattern of data found for contradicted belief-sets involving familiar definitions (See Elio, 1996). Cummins (1995) also noted that the few-disablers causal relations are quasi-definitional in nature.

It would be parsimonious to extend the notion of disablers (and hence the possible-worlds account) to the treatment of contradicted promises. For example, our understanding of a promise "*If you do x, then you will get y*" might include an understanding that many factors outside the control of the promiser might derail the promise. Thus, promises might be subsumed under a "many-disabling factors" category of relations. There is some support for this idea from the present data: the degree-of-belief ratings for promise-conditionals and promise-categorical statements are most similar to those given to causal-conditionals in the many-disabler cases, at least for the MP case. The prior degree-of-belief ratings also support this.

## Knowledge Types in Unfamiliar Domains

Does the influence of knowledge-types on belief revision require extensive domain knowledge? Or is there some domain-independent appreciation for causal relationships v. (merely) predictive relationships that may impact belief-revision decisions? In previous studies, Elio and Pelletier (1994; in press) used belief-revision problems with science-fiction cover stories and found that subjects were more likely to abandon the conditional on science-fiction problems than on the same problems using nonsense syllables. The study presented next offers a more deliberate investigation of the hypothesis that a reasoner's understanding of *type* of knowledge expressed in a conditional can influence how contradictions are resolved, even for unfamiliar domains.

### The conditionals

A set of fifteen conditional statements using anthropological concepts was constructed. The same subjects who categorized the conditionals used in the study described above also categorized the anthro conditionals as expressing a definition (the familiar v. unfamiliar distinction was dropped), a promise, cause-and-effect, and prediction. Of the fifteen conditionals, none were clearly identified as promises and four had no clear majority votes for any category. The remaining eleven were grouped into three types: four causal statements, four statements that were classed nearly equally often as predictions and definitions, and three statements that were clearly classed as predictions. While these classifications were not as consistent as I had hoped, this small set of conditionals allowed a preliminary study of how different knowledge-types might impact belief revision decisions for unfamiliar domains. Example anthro conditionals (with the number of classifications as <u>c</u>ausal, <u>p</u>rediction, <u>d</u>efinition, or <u>prom</u>ise) are: (a) *Causal:* If there is a death in a Meorian tribe, then the tribe relocates its camp. (c=25; p=10); (b) *Prediction/Definition:* If different families speak the same dialect of S'wara, then they belong to the same sharing camp. (p=15; d=21); (c) *Prediction:* If two villages are meeting for food exchanges, then the hosting village is represented by a female. (p=22; d=9; prom=9). The full set is available in Elio (1996).

### Subjects and Design

Thirty-nine subjects solved eleven problems based on the entire set of anthro conditionals. Because of the awkward number of conditionals of each type, the design of this preliminary study was not ideal: in each problem set, inference rule (MT v. MP) and knowledge type were repeated factors on just the causal and predict/define conditionals. The pure prediction conditionals appeared *either* in modus ponens or in modus tollens form to a particular subject.

**Results**

Table 3 gives the mean belief ratings for the anthro conditional and data statements as a function of knowledge type. First consider only the causal and prediction/definition problems. The findings here are consistent with what Elio and Pelletier previously found for an unfamiliar domain, namely science-fiction problems: (a) belief in the conditional was lower than belief in the data statement and (b) this was more pronounced for MP belief sets than for MT belief sets: On modus-ponens belief sets, subjects accorded conditionals a relatively lower belief than the data statements for the prediction/definitional scenarios (-2.15) than they did when the conditionals involved causal relationships (-.86).

The new result is that both these significant main effects interacted with knowledge type, indicating that, even for relatively unfamiliar domains, contradictions involving what were classified as causal relationships are reasoned about differently than those classified as non-causal relationships ($F(1,38)=6.79$, $p=.013$). The Table 3 results are, understandably, different than what was found for the familiar topics (Table 2), for which the reasoner has more information about the existence and relevance of other factors that may influence the subjective likelihood associated with a $p—>q$ statement.

| Modus Ponens Belief Set | | | |
|---|---|---|---|
| Upon learning $\sim q$, the new degree of belief in ... | | | |
| | $p —>q$ | $p$ | Difference |
| Causal | 3.9 | 4.7 | -.8 |
| Prediction/Definition | 3.2 | 5.3 | -2.1 |
| Prediction | 3.5 | 5.7 | -2.2 |

| Modus Tollens Belief Set | | | |
|---|---|---|---|
| Upon learning $p$, the new degree of belief in... | | | |
| | $p—>q$ | $\sim q$ | Difference |
| Causal | 4.2 | 4.7 | -.5 |
| Prediction/Definition | 4.2 | 4.0 | .2 |
| Prediction | 3.2 | 5.0 | -1.8 |

Table 3: Belief Ratings—Anthropology Problems

A separate analysis was done on the belief ratings given on the "pure prediction" problems (line 3 in Table 3), in which inference rule (MP v. MT) was a between-subjects factor. Following contradiction, these pure prediction conditionals received a significantly lower degree-of-belief rating than did the data statements, and this was more pronounced for modus ponens problems than for modus tollens problems ($F(1,37)=26.85$, $p < .001$). Comparing the pattern of results for the pure-prediction conditionals with the results for the causal conditionals, it again seems that even for unfamiliar domains, the relationship understood to hold between the antecedent and the consequent influences whether the conditional or the categorical belief is called into question, given a contradiction to a valid inference.

It can be argued that these anthro conditionals are not all that unfamiliar: the domain concerned human behaviors and customs, and people have enough general knowledge about such a domain to be able to conjecture whether a relation in this domain was causal or prediction. But that is fine. The understanding of a "merely" predictive relationship may include the knowledge that there are many possible worlds in which consequent does not co-occur with the antecedent; the reasoner may then consider it is possible that the current world is one such world, as way to reconcile the contradiction. What is important, at this point at least, is that people appear to be able to interpret sentences as signifying one *kind* of relationship between the antecedent and consequent, even when the domain is unfamiliar, and that distinction influences how a contradiction to a valid inference impacts the belief set.

## Summary

This task, the results reported here, and the possible-worlds interpretation relate to and extend several themes that have emerged in the deductive reasoning literature. The first concerns the manner in which the belief state is modeled— Elio & Pelletier (in press) discuss the parallels between the syntactic v. model-theoretic perspectives of belief states in the AI community with the mental logic v. mental model approach to deduction in the psychological perspective. The decision to adopt one approach or the other has implications for how one can define and measure formal notions like epistemic entrenchment or minimal change.

The second theme is the kind of inferences a reasoner finds plausible to make, since this determines the content of the belief state, even before contradictory information arrives. What is particularly relevant here is, on the one hand, data on the so-called "suppression" of valid inferences (e.g., Byrne, 1989) and, on the other hand, the surrounding theoretical considerations some researchers have put forward on this topic on this data. For example, Stevenson and Over (1995) discuss a role for weighted mental models that reflect a reasoner's assessment of the plausibility of alternative possible worlds. Earlier, Markovits (1984) argued that a reasoner's ability to generate possible alternative causes of an effect mediated the ability to reason correctly (generate only valid inferences) given deductive reasoning problems with conditional premises. The belief revision tasks, by presenting a contradiction, may prime a subject's consideration of possible worlds in which additional factors come into play. More generally, belief-revision and belief-update decisions (where "update" is an extension of the belief set when new information does not contradict existing beliefs) may invite a kind "belief-based" reasoning (George, 1995), in which subjects are not assuming that the given premises are true. The results presented here support the view that belief revisions are a function of knowledge about the relation between the antecedent and consequent. Thus, the many v. few disabler notion—as well as the "type of knowledge expressed"—may reduce to an understanding of the frequency with the consequent is entailed by the antecedent, in all the (plausible) worlds a reasoner can generate. This view is

consistent with the expansion of mental-models theories to include some aspect of uncertainty or subjective probability (e.g., Stevenson & Over, 1995; Johnson-Laird, 1994).

The final matter here concerns the difference in how belief-revision decisions are made when the initial belief state was defined with a modus ponens or a modus tollens inference. For familiar content (i.e., the data in Table 1), contradicted MP and MT belief sets were, by and large, revised somewhat differently. Crudely put, subjects considering contradictions to MP belief sets disbelieved the data statement; subjects considering contradictions to MT belief sets disbelieved the conditional statement. (The anomaly is the uniform abandonment of promise conditionals, regardless of the form of the initial belief set). Elio & Pelletier (in press) note that the MP belief-revision problem becomes something like an MT inference problem, and vice versa. From that perspective, it is curious that subjects do not merely opt for the "easy" MP inference afforded by the new information ($p$) with the old conditional ($p—>q$) to conclude q and hence disbelieve the old categorical premise ~q. That lead Elio and Pelletier to argue that this paradigm does not reduce to a deductive problem consisting of the old conditional and the new categorical statement. That said, it is possible that the relative subjective probability associated the conditionals in the initial belief sets is influenced by the categorical premise that co-occurred with them. Stevenson and Over (1995) speculate that, for a given $p—>q$ and $~q$ reasoning context, confidence that $p$ is true might be sufficiently high, such that it is more plausible to lower the belief in the conditional than to draw the modus tollens inference. Even though the MT inference was *provided* for the subjects in these belief revision tasks, there is this possibility that subjects had a lower initial belief in the conditional in the context of $~q$, before the contradictory information was even presented. The prior belief ratings already collected on the conditionals need to be augmented with belief ratings assigned in the context of the categorical premise. That said, the influence of alternative disablers and general knowledge-types on the belief revision decision for modus ponens and modus tollens problem sets still stands.

In summary, the belief revision "problem," (defined here as choosing one of possibly several previously-held beliefs to abandon or at least call into question, when faced with contradiction) crystallizes several key issues of everyday plausible reasoning. For normative models, epistemic entrenchment seems better modeled as a function of the plausibility of the alternative worlds the reasoner generates to account for a contradiction. For descriptive models of human reasoning, the "background knowledge" or "interpretative procedures" that are typically referenced in theoretical accounts of deductive reasoning take a central role modeling how contradictions to valid inferences are resolved.

## Acknowledgments

## References

Alchourrón, C., Gärdenfors, P., & Makinson, D. (1985). On the Logic of Theory Change. *Journal of Symbolic Logic,* 50, 510-530.

Byrne, R. (1989). Suppressing valid inferences with conditionals. *Cognition*, 61-83.

Cummins, D. D. (1995). Naive theories and causal deduction, *Memory & Cognition*, 23, 646-658.

Cummins, D.D., Lubart, T., Alksnis, O., & Rist, R. (1991). Conditional reasoning and causation. *Memory & Cognition, 19*, 274-282.

Elio, R. (1996). On the epistemic entrenchment of different knowledge types in belief revision tasks. (Tech. Rep. 96-16). Edmonton, Alberta: University of Alberta, Department of Computing Science.

Elio, R. & F. J. Pelletier (1994). The Effect of Syntactic Form on Simple Belief Revisions and Updates. *Proceedings of the 16th Annual Conference of the Cognitive Science Society.* Atlanta (pp. 260-265).

Elio, R. & Pelletier, F.J. (in press). Belief revision as propositional update. *Cognitive Science.*

Evans, J. St.T., Newstead, S. E., & Byrne, R. M. J. (1993). *Human reasoning.* Hillsdale, NJ: Lawrence Erlbaum.

Gärdenfors, P. (1988). *Knowledge in flux: Modeling the dynamics of epistemic states.* Cambridge, MA: MIT.

George, C. (1995). The endorsement of the premises: Assumption-based or belief-based reasoning. *British Journal of Psychology, 86*, 93-111.

Grice, P. (1975). Logic and conversation. In P. Cole & J. L. Morgan (Eds.), *Studies in syntax, Vol. 3: Speech Acts.* New York: Academic Press.

Foo, N.Y., & Rao, A.S. (1988) *Belief revision in a microworld* (Tech. Rep. No. 325). Sydney: University of Sidney, Basser Department of Computer Science.

Johnson-Laird, P. N. (1994). Mental models and probabilistic thinking. *Cognition,* 50, 189-209.

Markovits, H. (1984). Awareness of the 'possible' as a mediator of formal thinking in conditional reasoning problems. *British Journal of Psychology, 75*, 367-376.

Nebel, B. (1991). Belief revision and default reasoning: Syntax-based approaches. In *Proceedings of the Second Conference on Knowledge Representation,* 417-428, San Mateo, CA: Morgan Kaufmann.

Stevenson, R.J. & Over, D. E. (1995). Deduction from uncertain premises. *Quarterly Journal of Experimental Psychology, 484*, 613-643.

Thagard, P. (1989). Explanatory coherence. *Behavioral and Brain Sciences, 12* 435-502.