

## A Plausibility-Based Approach to Incremental Inference

**David J. Stracuzzi**

Sandia National Laboratories  
Albuquerque, NM 87185-1188

### Abstract

Inference techniques play a central role in many cognitive systems. They transform low-level observations of the environment into high-level, actionable knowledge which then gets used by mechanisms that drive action, problem-solving, and learning. This paper presents an initial effort at combining results from AI and psychology into a pragmatic and scalable computational reasoning system. Our approach combines a numeric notion of plausibility with first-order logic to produce an incremental inference engine that is guided by heuristics derived from the psychological literature. We illustrate core ideas with detailed examples and discuss the advantages of the approach with respect to cognitive systems.

### Introduction

Reasoning plays a fundamental role in cognition and intelligence. It supports decision making, action, problem-solving, and learning by converting observations of the surrounding environment into a high-level representation of the current situation. Cognitive systems aimed at large and complex environments therefore require the support of a computationally efficient and scalable inference system to provide these interpretations.

The artificial intelligence and psychology research communities have both studied reasoning extensively. Artificial intelligence has focused on developing both powerful logic-based and uncertainty-based computational mechanisms for reasoning under a variety of circumstances. However, the resulting systems typically rely on proving or on carefully estimating probabilities for large numbers of beliefs at great computational expense. The resulting systems tend to ignore plausible but unsound conclusions, to process large batches of data instead of individual observations, and often do not scale well. Conversely, psychologists tend to focus on describing at a high level the mechanisms that humans use for reasoning under specific circumstances. Most mechanisms therefore lack sufficient description a computational implementation, or are too specific to apply to a general cognitive systems.

This paper presents initial work on the Plausible Logic Inference Engine (PLIE), which seeks to combine results

from both artificial intelligence and psychology into an incremental, pragmatic, and scalable computational reasoning system. Our approach includes three key elements. First, it combines first-order logic with uncertainty based on a notion of plausibility to provide the system with a flexible knowledge representation. Second, the system includes an inference mechanism that integrates deductive and abductive methods to provide a more robust inference capability than either method can provide alone. Importantly, the inference mechanism also supports incremental update of individual beliefs. Finally, PLIE includes a guidance mechanism based on biases identified in the psychology literature to determine which inferences the system should pursue.

Our presentation focuses first on describing the knowledge representation and specific inference patterns applied by the inference engine. Next we discuss the inference process along with the heuristics used to guide it. We illustrate key ideas throughout using an expanded version of Pearl's (1988) burglar alarm example. We also discuss the relationship of our approach to cognitive systems in general throughout the presentation. Later, we outline the relationship of our proposed techniques to other inference methods and discuss next steps toward realizing the stated goals, including possible extensions to analogical reasoning, concept and structure learning, and metacognitive abilities such as bias learning.

### Knowledge Representation

Combining first-order logic with uncertainty has long been a goal of artificial intelligence. First-order logic provides a compact representation and efficient inference techniques, but it does not handle the uncertainties present in the real world. Probabilistic methods naturally capture and reason about this uncertainty, but they do not represent or exploit the relational structure of the world. Recent efforts at combining the two, such as Markov logic (Richardson and Domingos 2006), have shown success. However, such methods rely heavily on complex statistical computation which has not been shown to scale up.

This work uses a combination of logic and uncertainty based on plausibility as first described by Polya (1954), who codified the techniques used for hypothesizing theorems and guiding the search for the associated proofs. Polya's approach was based on a qualitative notion of accumulating

confidence from evidence, and is largely consistent with Dempster-Shafer theory (Shafer 1976), a mathematical theory of evidence that allows one to combine many sources of evidence to arrive at a degree of belief. Our approach replaces the qualitative account of confidence with a quantitative one, but retains the important properties highlighted by Polya. The result is similar to MYCIN (Shortliffe and Buchanan 1975) and to Friedman's (1981) work on plausible inference.

### Beliefs and Working Memory

PLIE represents knowledge using a combination of predicate logic and uncertainty. A *predicate* represents a generalized concept, or equivalently, a class of environmental situations. Each predicate may include a list of arguments. For example, Alarm(x) represents the situation in which the burglar alarm at location x is sounding.

*Beliefs* represent specific groundings or instances of a predicate such that each variable in the predicate's argument list is bound to some domain constant. For example, Alarm(BOB) indicates a belief that the alarm at Bob's house has sounded. Observations are similar to beliefs in that they correspond to grounded predicates. Whereas the inference process constructs beliefs, observations come from the perceptual process. This work does not speculate as to the nature of perception; we assume that some process produces the needed results.

Each belief or observation,  $b$  also has an associated plausibility score,  $-1 \leq \text{Pl}(b) \leq 1$ , such that  $\text{Pl}(b) = 1$  indicates a strong belief in  $b$  while  $\text{Pl}(b) = -1$  indicates a strong belief against  $b$ .  $\text{Pl}(b) = 0$  indicates no evidence in either direction and represents the default state of all beliefs. Plausibility scores represent the difference in evidence for and against the belief,  $\text{Pl}(b) = E^+(b) - E^-(b)$ , where  $0 \leq E^+(b), E^-(b) \leq 1$ . Importantly, the plausibility of a belief is tied to the plausibility of its logical negation,  $\text{Pl}(b) = -\text{Pl}(\neg b)$ , preventing contradictions. The methods for calculating these evidence values are described later.

PLIE stores both beliefs and observations in its *working memory* as nodes in a directed graph. The edges of the graph represent derivational information, and stem from the inference process. Specifically, a directed edge from belief  $b_1$  to  $b_2$  indicates that  $b_1$  formed part of the evidence used to conclude (or explain)  $b_2$ . Thus, observations have only outward directed edges, while all inferred beliefs have at least one edge directed inward.

### Rules and Long-Term Memory

Graph edges derive from the *rules* stored in *long-term memory*. Rules represent the relationships among predicates that can be used to derive new beliefs. From this perspective, the contents of working memory may be viewed as activated structures from long-term memory. This is consistent with Cowan's (1988) view of working memory, which holds that it is an extension of long-term memory, and that the number of activated long-term structures is not limited, although the number of items in the focus of attention is limited.

Rules in PLIE take the form

$$p_1(\cdot), p_2(\cdot), \dots, p_m(\cdot) \Rightarrow q_1(\cdot), q_2(\cdot), \dots, q_n(\cdot),$$

where  $p_i(\cdot)$  and  $q_j(\cdot)$  represent predicates, and  $\Rightarrow$  denotes *threshold-implication* (Stracuzzi and Könik 2008). Threshold-implication defines the consequent to be true if the linear-threshold function,  $\sum_{i=0}^m v_i \text{Pl}(p_i) > 0$ , is satisfied given that  $v_i$  represents the weights associated with each term in the antecedent ( $w_i$  for the consequent), and  $\text{Pl}(p_0) = \text{Pl}(q_0) = 1$  by definition. Note that  $p_i(\cdot)$  and  $q_i(\cdot)$  represent general predicates while  $p_i$  and  $q_i$  denote the specific groundings used to instantiate a rule. which PLIE identifies by pattern matching from working memory. In this work, both sides of the implication represent linear threshold functions, so when the antecedent function is satisfied, the system adds or updates beliefs as needed to satisfy the consequent function.

Each implication also has two associated parameters that represent rule strength. For a rule  $A \Rightarrow B$ , the parameter  $\alpha \approx \text{Pr}(B|A)$  represents the reliability of the rule when viewed deductively, meaning that belief in A implies a belief in B. Similarly,  $\beta \approx \text{Pr}(A|B)$  represents the reliability of the rule when viewed abductively, meaning that belief in B allows the assumption of A. Here,  $0 \leq \alpha, \beta \leq 1$ , though the two are not complimentary. In practice these parameters will initially take default values, which the system can later modify based on experience.

Table 1 shows the rules for the burglar alarm example. For the sake of clarity, they are expressed using Boolean connectives rather than in threshold logic. For example, PLIE could represent the antecedent of rule (6) as

$$1 \cdot \text{Pl}(\text{SecurityCo}(y, x)) + 1 \cdot \text{Pl}(\text{Calls}(y, x)) - 1.5 > 0.$$

Threshold logic provides several important advantages. First, it supports a more powerful representation than Horn clauses, including conjunctions, disjunctions, and many other simple functions such as  $p_1 \wedge (p_2 \vee p_3)$ . The underlying numeric representation (the weights) also provide a natural way to handle numeric (as opposed to symbolic) domain constants. Finally, with respect to the long-term goals of this work, Stracuzzi and Könik (2008) demonstrate how

- 
- |     |   |
|-----|---|
| (1) | Burglar( $x_1$ ) $\xRightarrow[\beta=0.40]{\alpha=0.95}$ Alarm( $x_1$ )   |
| (2) | Earthquake( $x_2$ ) $\xRightarrow[\beta=0.60]{\alpha=0.30}$ Alarm( $x_2$ )  |
| (3) | Alarm( $x_3$ ) $\wedge$ Neighbor( $y_3, x_3$ ) $\wedge$ Quiet( $y_3$ )<br>$\xRightarrow[\beta=0.05]{\alpha=0.90}$ Calls( $y_3, x_3$ )       |
| (4) | Alarm( $x_4$ ) $\wedge$ Neighbor( $y_4, x_4$ ) $\wedge$ $\neg$ Quiet( $y_4$ )<br>$\xRightarrow[\beta=0.5]{\alpha=0.70}$ Calls( $y_4, x_4$ ) |
| (5) | PhoneRings( $x_5$ ) $\wedge$ Neighbor( $y_5, x_5$ ) $\wedge$ Quiet( $y_5$ )<br>$\xRightarrow[\beta=0.95]{\alpha=0.10}$ Calls( $y_5, x_5$ )  |
| (6) | SecurityCo( $y_6, x_6$ ) $\wedge$ Calls( $y_6, x_6$ ) $\xRightarrow[\beta=1.0]{\alpha=1.0}$ Burglar( $x_6$ )                                |
- 

Table 1: Long-term memory contents for the alarm example.

<b>Pattern</b>	<b>Symbolic Form</b>	<b>Evidence Update Equation</b>
Modus Ponens • <i>Deductive</i> • <i>Increases evidence for <math>q_i</math></i>	$p_i(\cdot) \Rightarrow q_j(\cdot)$ <hr/> $\sum_{i=0}^m v_i \text{Pl}(p_i) > 0$ <hr/> $\sum_{i=0}^n w_i \text{Pl}(q_i) > 0$	$\Delta E^+(q_i) = \alpha \frac{\sum_{j \in S} v_j \text{Pl}(p_j)}{\sum_{j \in S}  v_j }$
Modus Tollens • <i>Deductive</i> • <i>Increases evidence against <math>p_i</math></i>	$p_i(\cdot) \Rightarrow q_j(\cdot)$ <hr/> $\sum_{i=0}^n w_i \text{Pl}(q_i) < 0$ <hr/> $\sum_{i=0}^m v_i \text{Pl}(p_i) < 0$	$\Delta E^-(p_i) = -\alpha \frac{\sum_{j \in S} w_j \text{Pl}(q_j)}{\sum_{j \in S}  w_j }$
Denying the Antecedent • <i>Deductive</i> • <i>Increases evidence against <math>q_i</math></i>	$p_i(\cdot) \Rightarrow q_j(\cdot)$ <hr/> $\sum_{i=0}^m v_i \text{Pl}(p_i) < 0$ <hr/> $\sum_{i=0}^n w_i \text{Pl}(q_i) < 0$	$\Delta E^-(q_i) = -\alpha \frac{\sum_{j \in S} v_j \text{Pl}(p_j)}{\sum_{j \in S}  v_j }$
Accepting the Consequent • <i>Abductive</i> • <i>Increases evidence for <math>p_i</math></i>	$p_i(\cdot) \Rightarrow q_j(\cdot)$ <hr/> $\sum_{i=0}^n w_i \text{Pl}(q_i) > 0$ <hr/> $\sum_{i=0}^m v_i \text{Pl}(p_i) > 0$	$\Delta E^+(p_i) = \beta \frac{\sum_{j \in S} w_j \text{Pl}(q_j)}{\sum_{j \in S}  w_j }$

Table 2: Inference patterns and associated confidence updates. For the symbolic forms, note that the implication (first line) represents a rule from long-term memory, the inequality above the horizontal bar represents the evidence required for the pattern to apply, and the inequality below the bar represents the conclusion drawn given the pattern, the rule, and the evidence.

to learn a set of hierarchically structured predicates based on threshold logic functions from sparse examples.

## Inference Patterns and Confidence Propagation

Researchers in artificial intelligence have considered inference methods extensively, though often in a piecemeal fashion. As noted earlier, the integration of logical and statistical inference has only recently come to the fore with methods such as Markov Logic (Richardson and Domingos 2006) and Bayesian Logic Programs (Kersting and De Raedt 2005). The integration of multiple inference patterns, such as deduction, abduction, and analogy has received even less attention. For example, the Markov logic framework includes a deductive inference engine which others have demonstrated can be used for abduction (Kate and Mooney 2009). However, the two methods rely on different rule structures so that the user must decide which method to apply before encoding the domain rules.

In this paper, we consider the combination of deduction and abduction. The utility in combining the two follows from their complementary nature. Deduction produces sound arguments (in which the conclusion necessarily holds given the evidence), but can only restate existing knowledge in different terms. Abduction adds new information to the system by hypothesizing (unsound) explanations for existing knowledge, which deductive methods can then further expand upon. The resulting system should therefore reason more broadly than one based on either method alone.

## Inference Patterns

PLIE relies on four distinct inference patterns for constructing new beliefs. Each pattern, shown in Table 2, applies to specific situations depending on the available evidence. For example, *modus ponens* deductively concludes the consequents ( $q_j$ ) of a rule given that the antecedents ( $p_i$ ) hold. In the context of threshold-implication, this means that the evidence equation,  $\sum_{i=0}^m v_i \text{Pl}(p_i) > 0$ , must be satisfied by pattern matching beliefs from working memory to the predicates in the rule antecedent. When successful, PLIE concludes that  $\sum_{i=0}^n w_i \text{Pl}(q_i) > 0$  must also hold. In practice, this means that the system increases the plausibility in each of the rule's consequent beliefs,  $q_i$ . The other inference patterns operate analogously depending on which beliefs get used as evidence.

The four inference patterns have differing characteristics. Two of them, modus ponens and modus tollens, are logically sound, deductive rules. A third, denying the antecedent, is also deductive in nature though not logically sound because it ignores the possibility that a conclusion may be drawn in more than one way. Nevertheless, logically unsound inference patterns are still useful for accumulating evidence about beliefs (Polya 1954) and for considering beliefs that, while not provable, may still accurately describe an agent's environment. Finally, affirming the consequent serves as the basis for abduction, which is also not logically sound for similar reasons. Notice also that two of the patterns deal with inference from positive evidence (plausibilities sufficiently greater than zero), while two of them deal with negated evidence (plausibilities less than zero). This lets PLIE handle a broad array of situations and rule structures.

Consider the following application of the inference pat-

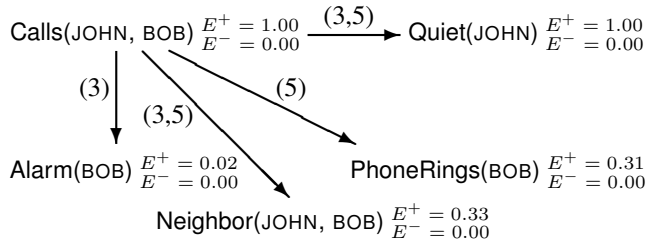


Figure 1: Contents of working memory generated by applying the inference patterns and evidence updates from Table 2 to the long-term memory contents in Table 1. Edge labels indicate the rules used to derive the beliefs.

terns to the rules shown in Table 1. For now, we consider only the symbolic application of patterns, but return in the next section to describe the plausibility updates. Suppose that working memory initially contains the observations  $Quiet(John)$  and  $Calls(John, Bob)$ , each with plausibilities of 1.0. The system can then apply rules (3) and (5) with accepting the consequent to conclude (or assume)  $Alarm(Bob)$  (rule 3),  $PhoneRings(Bob)$  (rule 5), and  $Neighbor(John, Bob)$  (both rules). Figure 1 shows the coherence graph that would result if the system applied both rules. Notice that  $Calls(John, Bob)$  explains  $Quiet(John)$ . The opposite is not true in this case because the antecedents to rules (3) and (5) are not satisfied, which prevents PLIE from applying the rules with modus ponens.

Similarly, rule (4) does not apply in this case. Although the consequent matches with the observation  $Calls(John, Bob)$  in working memory, the antecedent term  $\neg Quiet(John)$  contradicts the observation  $Quiet(John)$ . For now, we assume for simplicity that PLIE cannot change the plausibility of observations (it trusts its senses). However, the assumption is not fundamental, and may be dropped in later versions of the work. Rule (6) also does not apply because, in the context of working memory contents, it fits the modus ponens pattern, which requires that the antecedent be satisfied.

## Evidence Updates

As discussed earlier, plausibility quantifies the system's belief in specific aspects of the environment. Specifically, PLIE represents and propagates any uncertainty associated with its rules and observations through the inference process. This is particularly important in the context of abduction, which produces logically unsound assumptions. Some of these may be reliable, while others may not. Quantifying this uncertainty lets PLIE distinguish between these two cases and focus on working with more plausible conclusions during subsequent inference steps. This is an essential aspect of maintaining tractability in complex domains.

The update equations shown in Table 2 implement this by first combining the plausibility values used to trigger the rule, and then dividing it among the conclusions. The summation factor in the equations determines the amount of evidence that gets propagated by performing a weighted sum over the plausibility scores associated with the matched be-

liefs. However, recall that threshold logic can represent disjunctive (or partially disjunctive (as in  $p_1 \wedge (p_2 \vee p_3)$ ) rule bodies. This means that not all of the terms in the evidence must be satisfied for the rule to apply.

A belief  $b$  satisfies a rule term if  $\text{sign}(w_b \text{Pl}(b))$  is positive for the inference patterns that increase confidence in conclusions ( $b$  contributes to exceeding the threshold), and negative for patterns that decrease confidence in conclusions ( $b$  contributes to *not* exceeding the threshold). We define  $S$  as the set of beliefs used to instantiate the *satisfied* terms in the evidence of a rule application. For example, suppose  $p_1 \wedge (p_2 \vee p_3)$  represents the evidence for an inference step and that PLIE has  $p_1$  and  $p_3$  in working memory with positive confidence, and  $p_2$  with negative confidence. Applying modus ponens would produce  $S = \{p_1, p_3\}$ , while denying the antecedent would produce  $S = \{p_2\}$  (and fail).

Before PLIE propagates the accumulated update value to the conclusions, the value first gets scaled by the reliability of the applied rule ( $\alpha$  or  $\beta$ ) as shown in Table 2. This lends greater confidence to conclusions drawn from reliable rules. Note also the sign of the evidence update, which is negated for inference patterns that increase evidence against a belief. In these cases, the weighted sum of the evidence values will be negative because the sign of the weighted plausibilities is negative. Evidence values are always positive however, so the sign needs to be flipped.

Finally, the evidence score for a belief  $q_i$  is updated by

$$E_{t+1}(q_i) = (1 - \lambda \hat{w}_i) E_t(q_i) + \lambda \hat{w}_i \Delta E_t(q_i) \quad (1)$$

where  $\hat{w}_i = \frac{|w_i|}{\|w\|}$  represents the normalized weight magnitude associated with  $q_i$  in the rule, and  $0 < \lambda < 1$  represents a primacy-recency trade-off parameter (discussed below). The equation for updating  $p_i$  is analogous. Note that we omitted the positive and negative superscripts from  $E$ ; Table 2 shows which evidence value gets updated by a given inference pattern. However, if the  $q_i$  term is negated in the rule ( $w_i < 0$ ), then the update gets switched from  $E^+$  to  $E^-$  or vice versa. This follows from the plausibility relationship  $\text{Pl}(b) = -\text{Pl}(\neg b)$  between beliefs and their negations.

Equation 1 implements a form of the primacy-recency bias (Ebbinghaus 1913) by scaling the existing and updated evidence scores. Setting  $\lambda$  near 1 causes new evidence to have greater impact on plausibility (recency effect), while  $\lambda$  near 0 minimizes the impact of new evidence (primacy effect). Psychologists view these as memory biases, but this work assumes that biases in how beliefs and their supporting evidence get stored impact the reasoning process. Along similar lines, rule terms with relatively large weight magnitudes receive higher impact from new evidence than those with low weight magnitude, as these contribute the most to satisfying the rule.

Returning now to Figure 1, notice that the plausibility associated with  $Alarm(Bob)$  is very low at 0.04 (assuming that  $\lambda = 1$  and the weight on each rule term is 1). This follows from the low  $\beta$  value associated with rule (3), implying that John has a high false-positive rate. The plausibility of  $PhoneRings(Bob)$  is much higher, indicating that John's call is a much better indicator of Bob's phone ringing than of his burglar alarm's sounding. Rules (3) and (5) both contribute

to the plausibility of Neighbor(JOHN, BOB) (assume for now that the rules were applied in the order listed in Figure 1).

Upon initial review, the inference patterns and evidence updates may appear complex. However, closer inspection reveals that all of the patterns and equations represent minor variations on a general mechanism depending on whether the inference flows with or against the rule's implication sign, and whether the evidence increases the evidence for or against a belief. This closely resembles Friedman's (1981) work, except that in switching from Boolean to threshold logic, we have collapsed several highly specialized cases into one general case. As a result, the implementation details should be substantially simpler.

### Inference Engine and Heuristics

The inference patterns described above govern the details of how PLIE combines beliefs with rules to derive new beliefs. Given these patterns, the inference engine first selects beliefs and rules from memory, and then applies the patterns to construct new beliefs and update plausibility scores. Importantly, the incremental application of the inference patterns from Table 2 provides a seamless integration of forward and backward reasoning methods, allowing PLIE to exploit opportunities in reasoning as they become available. A similar idea appears in the planning community as opportunistic planning, or metaplanning (see work by Hayes-Roth and Hayes-Roth 1979 for an early example). In this section, we describe the details of the inference engine, including the methods used to select beliefs and rules for inference.

### Heuristics and Biases

This work aims to produce a scalable inference engine that, while neither complete nor sound, performs well in practice and provides a foundation for other cognitive abilities. At present, PLIE does not have access to an agent's goals, so our approach focuses on properties of beliefs based on three cognitive biases that specifically impact reasoning and for which substantial evidence has been amassed. We do not claim to reproduce the exact conditions that lead to these biases as identified in the literature. Instead, we focus on incorporating a generalized view of each bias into our computational reasoning system. Incorporating other biases and heuristics remains a key area of future work.

One such property, explanatory coherence (Thagard 1989), implies a preference for beliefs that "hang together" in the context of explaining observations. Although often associated with abduction, we also apply coherence to deduction, as PLIE may also make unsound deductive inferences.

PLIE implements coherence as a measure of linkage with emphasis on the plausibility of the linked beliefs and their proximity to observations. Let  $\text{adj}(b)$  represent the set of beliefs adjacent to  $b$  in working memory regardless of edge direction. For a belief  $b' \in \text{adj}(b)$ , let  $\text{rel}(b', b)$  represent the rule reliability ( $\alpha$  or  $\beta$ ) that applies when the system views the rule that links  $b$  with  $b'$  as though it were used to derive  $b$  from  $b'$  (see Table 2). Also let  $d(b')$  represent the shortest distance from  $b'$  to any observation in working memory (0 if  $b'$  is an observation).

Given these definitions, we define coherence as

$$\text{coh}(b) = \sum_{b' \in \text{adj}(b)} \frac{\hat{w}_{b'} \text{rel}(b', b)}{d(b') + 1} \text{Pl}(b').$$

Though related to plausibility in its use of belief weights and rule reliabilities, coherence is a measure of the total evidence in working memory for  $b$ , with each piece of evidence weighted by its proximity to an observation. In practice, a unit increment in the coherence measure represents direct evidence from or explanation of one (high plausibility) observation through a reliable rule. The measure attenuates as evidence gets divided over multiple rule terms, rules become unreliable, observations grow more distant, and evidence terms become less plausible. As a result, coherence is unbounded in magnitude and cannot be computed incrementally. Nevertheless, the coherence measure is computationally simple and calculated locally, so the approach should remain tractable as the size of working memory grows large. Our coherence measure is similar in spirit to that of Ng and Mooney (1990), but operates over individual beliefs rather than proof trees.

The second bias implements a preference for working with highly plausible beliefs. We view this as a generalization of several specific biases identified in the literature, such as the confirmation and disconfirmation biases (Lord, Ross, and Lepper 1979). The former states that people tend to favor information that confirms their preconceptions, while the latter states that people tend to be very critical of information that contradicts them.

In practice, PLIE implements this in two ways. First, the system ignores beliefs that have low plausibility. In the alarm example, an agent may ignore the contingency in which the alarm has sounded if John's unreliability causes a sufficiently low plausibility score and no other evidence is available to increase it. Second, PLIE prefers rule applications with stronger evidence over those with weaker evidence. The extent to which a rule application exceeds its threshold indicates how strongly the evidence supports the rule application. This favors rule instances with more matched evidence terms, and whose evidence terms have high plausibility.

Finally, primacy and recency biases have already been implemented by the plausibility update equations. However, we apply a second form of recency in the inference engine by preferring to expand upon beliefs and observations that have been recently updated or added to working memory. This gets implemented by having the inference engine select only a small number of beliefs for expansion during any given inference cycle based largely on their the recency.

This combination of biases should cause PLIE to reason along a small number of plausible trajectories. Nothing explicitly prevents the system from exploring other lines of inference opened by new observations or substantial changes to the plausibility of established beliefs. However, it should not explore broadly through highly implausible regions of the belief-space.

### Inference Procedures

PLIE operates in cycles, with each cycle consisting of three steps. First, the system selects a small subset of beliefs from working memory for expansion and identifies the relevant rules from long-term memory. Next, it applies the inference patterns from Table 2 to derive new and update existing beliefs. Finally, it filters the resulting beliefs to remove those with low plausibility or coherence. The remainder of this section details each step in turn.

Each cycle begins with PLIE selecting a small subset of beliefs,  $W_{sel}$ , from working memory for expansion, where  $|W_{sel}|$  is an adjustable parameter. It selects beliefs with a preference for those most recently updated, breaking ties in favor of beliefs with higher coherence scores. PLIE then identifies the applicable rules from long-term memory given the beliefs in  $W_{sel}$ . A rule is applicable if three conditions are met: (1) a belief from  $W_{sel}$  matches a term in the rule evidence, (2) the remaining evidence terms can be matched by beliefs in working memory, and (3) the intended inference has not been previously made using the same rule and evidence. After identifying the set of applicable rules, PLIE then chooses the subset that exceeds their thresholds by the highest amounts such that each belief in  $W_{sel}$  gets paired with one applicable rule.

Given the selected beliefs and applicable rules, PLIE next applies the inference patterns from Table 2 to generate conclusions and update plausibility scores. This may entail multiple rule applications and may generate many conclusions, or multiple updates to a single conclusion. Though the search remains focused on recent, coherent beliefs, this provides a broader search through belief space than if the system used only the single “best” belief in each cycle.

Also note that conclusions drawn abductively by PLIE may contain unbound variables, such as  $Neighbor(y_3, BOB)$ . Here, the system creates a Skolem constant as a place holder. During subsequent inference cycles, it can unify Skolem constants with other constants. For example, this would effectively replace the belief “Bob has some specific, unnamed neighbor” with “Bob’s neighbor is John.” All instances of a Skolem constant must be replaced with the same domain constant and all unifications must be consistent (must not create a belief contradicts any existing belief).

After generating conclusions and updating plausibility scores, PLIE filters its results to remove beliefs with very low magnitude plausibility. Specifically, any belief  $b$  with  $|Pl(b)| < \tau$  does not get added to (or replaced in) working memory, where  $\tau$  is an adjustable parameter. This prevents it from creating and storing a large number of beliefs with scores very near to zero. Although PLIE is unlikely to select such beliefs for  $W_{sel}$ , they may still get matched as part of the inference process, thereby deriving a potentially large number of new beliefs with low plausibility.

An important side effect of this is that some beliefs may lose their support over time. The plausibility of such a belief does not change in working memory, but the coherence declines, which biases PLIE against selecting the belief for expansion. Plausibility changes only when the system draws a direct conclusion about a belief.

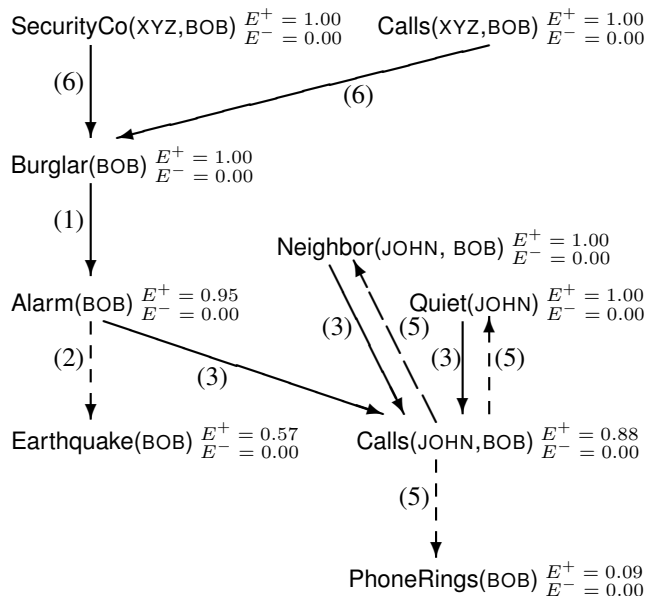


Figure 2: Coherence graph resulting from three inference cycles given the observations Neighbor(JOHN,BOB), Quiet(JOHN), SecurityCo(XYZ,BOB), and Calls(XYZ,BOB). Dashed arrows indicate available inferences not made by PLIE.

### An Example

Consider now a larger example based on the rules in Table 1. To begin, suppose that PLIE receives four observations, each with plausibility 1: Neighbor(JOHN,BOB), Quiet(JOHN), SecurityCo(XYZ,BOB), and Calls(XYZ,BOB). Initially, the only possible inference applies rule (6) with modus ponens to conclude that there is a burglar at Bob’s house, Burglar(BOB), as shown in the coherence graph in Figure 2. The high plausibility of the conclusion follows from the use of both highly plausible observations as evidence, and a highly reliable inference rule.

In the second cycle, PLIE chains forward from Burglar using rule (1) with modus ponens to conclude that the alarm has sounded at Bob’s house, Alarm(BOB). This conclusion is also highly plausible for reasons similar to those above. Note however that the coherence of Burglar (1.13) is substantially higher than that of Alarm (0.48) because Burglar links directly to observations.

PLIE has available two applicable rule instances in the third cycle given that  $W_{sel}$  contains Alarm(BOB). Rule (2) applies abductively and rule (3) applies with modus ponens. Assuming similar thresholds for the two rules, PLIE would prefer rule (3) over rule (2) based on the slightly larger plausibility contributions from Neighbor and Quiet (1.0 versus 0.95 from Alarm). PLIE therefore infers Calls(JOHN,BOB), which has a coherence of 1.16 based on its links to two observations. The inference also raises the coherence of Alarm slightly to 0.50. The small increment in this case follows from the low reliability of rule (3) when viewed abductively ( $\beta = 0.05$ ).

To summarize the inference process so far, PLIE has determined that a burglar has very likely struck Bob's house given that the XYZ security company has called him. His alarm system has also plausibly sounded, which implies that John may also call soon. PLIE avoided the inference that an earthquake triggered the alarm, though, even if made, the inference would be both less plausible and less coherent than the burglar explanation. Nevertheless, this does raise a question about whether preferences and coherence are sufficient to avoid or distinguish between multiple, conflicting explanations for a single set of observations. PLIE does not explicitly represent the notion that, while not technically mutually exclusive, rules (1) and (2) should not typically occur together. For now we retain this question as a point of future work.

A second open question relates to a stopping criterion. For cases in which the stream of incoming observations stagnates for a period of time, PLIE would ideally determine a point at which continuing inference is no longer productive. In the current example, PLIE could continue reasoning by applying rule (5) abductively to generate the assumption PhoneRings(BOB). As with Earthquake, this represents an alternate, and less plausible explanation for John calling Bob. If other rules related to events or beliefs stemming from the alarm's sounding or John's call to Bob were present in long-term memory, then the system would tend to follow those paths. In this case, no other inference paths are available, so PLIE returns to considering the low-plausibility and low-coherence lines of reasoning ignored earlier.

## Discussion

The plausibility-based approach to computational reasoning outlined above represents a distinct departure from traditional proof-based logic. It provides mechanisms for drawing unsound conclusions, relies on incomplete heuristic search. Although this admits the possibility that the system may draw incorrect conclusions or ignore many correct conclusions, it also represents an attempt at avoiding the intractability associated with many logical inference systems by identifying potentially useful conclusions. This is an important point in the context of cognitive agents whose goal is to perform tasks in complex environments.

From this perspective, adding information about agent goals to PLIE would be a fruitful direction for future work. For example, PLIE could use the relationship of a belief's predicate to goal predicates as an additional guidance heuristic. This respects the approach of exploring from observations (as well as goals) while further sharpening the system's focus on working with beliefs that are relevant to the task or domain at hand.

In the context of symbolic inference, Bridewell and Langley's (2011) AbRA is most similar to PLIE. Although their discussion focuses on abduction, AbRA can also apply rules deductively. Both systems can therefore combine both forward chaining from observations and backward chaining from goals. Their approach does not support inference from negated evidence (*modus tollens* and denying the antecedent) and does not handle uncertainty in beliefs or rules,

however. Friedman's (1981) work does support all four inference patterns identified in Table 2 and uses a numeric confidence measure that is somewhat similar to PLIE's plausibility scores, but relies on brute force computation to derive all possible beliefs.

Production systems such as Soar (Laird, Newell, and Rosenbloom 1987) also perform a kind of symbolic inference. Unlike most logical or even statistical inference systems, production systems do not commit to any one reasoning formalism, although the rules can be structured to implement a variety of formalisms. Instead, they take a purely syntactic approach, firing rules that match the current belief state and using other productions to implement preferences while typically ignoring the structural properties of the generated beliefs. In contrast, PLIE commits to a fixed reasoning formalism based on deduction and abduction while its guidance heuristics take the relationships among beliefs into consideration.

PLIE also represents a departure from traditional statistical optimization-based inference techniques. These methods, including Bayesian networks (Pearl 1985) and more recent efforts like Markov logic (Richardson and Domingos 2006), attempt to estimate probability distributions associated with one or more query variables. Statistical inference algorithms can require substantial computation, often drawing on information from large portions of the underlying belief network in response to a single query, or evaluating every possible belief simultaneously. Conversely, PLIE updates selected plausibilities incrementally in response to changes in observations. Though clearly less precise than statistical inference methods, PLIE entails a much lighter computational burden. This makes sense in the context of agents acting in worlds where responsiveness often takes precedence over precision.

The most immediate points of future work on PLIE concern implementation and performance comparison to other inference systems. This paper describes a first attempt at specifying a plausibility-based approach to a tractable and pragmatic reasoning system. The examples discussed above suggest that PLIE may demonstrate a number of interesting and useful properties, but a much more extensive evaluation is required. Of specific concern are the system parameters and the inference heuristics. With additional experimentation, the parameters should be set to fixed, domain independent values. Likewise, the impact of the inference heuristics requires further study. We selected plausibility, coherence, and primacy/recency because they are well-known and easy to implement. However, other implementations are possible and many other cognitive biases have been studied. The impact of these and other biases on system performance needs further study.

As noted earlier, the addition of analogical and inductive reasoning processes will play major roles in producing a system that scales and performs well on real world problems. Briefly, the role of analogy is to expand the applicability of the rules in long-term memory. By mapping known rules into novel domains, PLIE can leverage a greater proportion of its knowledge when it encounters new situations. Likewise, an inductive process will expand and add nuance to

the system's existing knowledge. This includes tuning rule strengths ( $\alpha$  and  $\beta$ ), modifying rule antecedents and consequents, and adding new rules based on experience. Together, these should substantially reduce PLIE's reliance on manual knowledge engineering and increase the breadth of its reasoning capability.

A third area of future work is metacognitive in nature. Currently, PLIE contains a set of four inference patterns, which may expand as analogical and inductive processes get added. These patterns are fixed, but could be adapted or expanded based on experience. For example, if the abductive pattern regularly produces conclusions that later get removed or refuted, then the system could modify the pattern require evidence with higher plausibility scores. This would increase its efficiency by finding and exploiting general patterns in the reasoning process.

### Concluding Remarks

The goal of this work is to combine ideas and results from a broad array of work from both the artificial intelligence and psychological research communities into a pragmatic and computationally tractable inference system. Our approach draws on experiences in both logical and statistical inference while dropping obsessions with formal proof and statistical precision. Likewise, even though our approach takes inspiration from research on human cognition, our objective is simply to identify reasoning heuristics that lead the inference engine to construct and maintain beliefs that are relevant to an agent acting in the environment.

The work reported here summarizes only an initial effort at designing a knowledge representation and inference algorithms to achieve the stated goals. The provided examples demonstrate the intended function of the system, but further testing is required to establish that the intended performance holds up across a variety of domains. Nevertheless, combination of logic, plausibility, and heuristics presented above provide strong foundation on which to experiment with, revise, and expand the core ideas of this work.

### Acknowledgments

Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.

### References

- Bridewell, W., and Langley, P. 2011. A computational account of everyday abductive inference. In *Proceedings of the Thirty-Third Annual Meeting of the Cognitive Science Society*. Boston, MA: Cognitive Science Society, Inc.
- Cowan, N. 1988. Evolving conceptions of memory storage, selective attention, and their mutual constraints within the human information-processing system. *Psychological Bulletin* 104(2):163–191.
- Ebbinghaus, H. 1913. *Memory: A Contribution to Experimental Psychology*. New York: Teachers College, Columbia University.
- Friedman, L. 1981. Extended plausible inference. In *Proceedings of the 7th International Joint Conference on Artificial Intelligence*, 487–495. Vancouver, BC, Canada: William Kaufmann.
- Hayes-Roth, B., and Hayes-Roth, F. 1979. A cognitive model of planning. *Cognitive Science* 30:275–310.
- Kate, R. J., and Mooney, R. J. 2009. Probabilistic abduction using markov logic networks. In *Proceedings of the IJCAI-09 Workshop on Plan, Activity and Intent Recognition*. Pasadena, CA: AAAI Press.
- Kersting, K., and De Raedt, L. 2005. Bayesian logic programming: Theory and tool. In Getoor, L., and Taskar, B., eds., *An Introduction to Statistical Relational Learning*. MIT Press. 291–321.
- Laird, J. E.; Newell, A.; and Rosenbloom, P. S. 1987. Soar: An architecture for general intelligence. *Artificial Intelligence* 33(1):1–64.
- Lord, C. G.; Ross, L.; and Lepper, M. R. 1979. Biased assimilation and attitude polarization: The effects of prior theories on subsequently considered evidence. *Journal of Personality and Social Psychology* 37(11):2098–2109.
- Ng, H. T., and Mooney, R. J. 1990. On the role of coherence in abductive explanation. In *Proceedings of the Eighth National Conference on Artificial Intelligence*, 337–342. Boston, MA: AAAI Press.
- Pearl, J. 1985. Bayesian networks: A model of self-activated memory for evidential reasoning. In *Proceedings of the Seventh Conference of the Cognitive Science Society*.
- Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Mateo, CA: Morgan-Kaufmann.
- Polya, G. 1954. *Mathematics and Plausible Reasoning, Volume II: Patterns of Plausible Inference*. Princeton, NJ: Princeton University Press.
- Richardson, M., and Domingos, P. 2006. Markov logic networks. *Journal of Machine Learning Research* 62(1–2):107–136.
- Shafer, G. 1976. *A Mathematical Theory of Evidence*. Princeton, NJ: Princeton University Press.
- Shortliffe, E. H., and Buchanan, B. G. 1975. A model of inexact reasoning in medicine. *Mathematical Biosciences* 23(3–4):351–379.
- Stracuzzi, D. J., and Könik, T. 2008. A statistical approach to incremental induction of first-order hierarchical knowledge bases. In Železný, F., and Lavrač, N., eds., *Proceedings of the Eighteenth International Conference on Inductive Logic Programming, LNAI 5194*, 279–296. Prague, Czech Republic: Springer-Verlag.
- Thagard, P. 1989. Explanatory coherence. *Behavioral and Brain Sciences* 12:435–467.