

A Coherence-Driven approach to Action Selection

Sindhu Joseph and Carles Sierra and Marco Schorlemmer

Artificial Intelligence Research Institute, IIIA-CSIC

Bellaterra (Barcelona), Catalonia, Spain

{joseph,sierra,marco}@iia.csic.es

Abstract

In this paper, we propose a coherence-driven approach to action selection in agents. The mechanism is inspired by the cognitive theory of coherence as proposed by Thagard. Based on a proposal to extend BDI agents with coherence, we interpret, how action selection can be viewed as a coherence-maximising problem. Contrasted against the classical BDI approach to action selection where actions are selected against a pre-determined set of beliefs and desires, this method offers a dynamic view of the cognitions of an agent, where a set of beliefs, desires and intentions are selected together to keep the coherence of the agent. We illustrate the approach by simulating how a coherence-driven robot selects its next action to pursue.

1 Introduction

A BDI-based reasoning process consists of a deliberative cycle in which an agent decides what state of affairs it wants to achieve from among all those desirable states of affairs [4; 15; 14]. The output of the deliberation process is a set of intentions (desires that the agent wants to pursue paired with a ‘top-level’ plan of action) [1]. Once the intentions are created and their associated preconditions (in the form of a set of beliefs) are met, then it is immediate that these intentions are realised.

As it should be apparent, there are a few major difficulties with this kind of reasoning. Among the many alternatives, it is not clear how a particular desire or a set of desires are chosen to be pursued further. In a graded cognitive agent, this could be done simply by selecting the desire that has the highest degree [3]. However, such a selection will not guarantee that the chosen desire is the best to be satisfied. To qualify for it, a desire should be consistent with most of the fundamental beliefs of the agent, and it should not conflict with other desires which are already in pursuit. Finally, this desire should be realisable. The last point is taken into account in the BDI deliberation cycle, however, not during the selection of the desire, rather at the point where intentions are generated [4; 15]. At this point, conflicts with other intentions may be discovered. In such cases, the plan is aborted and another plan or another desire itself has to be chosen. However, this is a

very ad-hoc procedure, and unlikely to result in an optimised or coherent agent.

Alternatively, our intuition says that among the many alternatives, a desire should be selected that is not only most desired, but also most coherent with the agent’s set of beliefs, other desires, and plans. The same is true to incorporate a new perception (belief). A new perception is incorporated only when it is coherent within the set of cognitions. The same happens when adopting a plan. That is, the model is essentially dynamic, where beliefs, desires and intentions are subject to the criterion of coherence maximisation. Here we propose to incorporate such a reasoning to artificial agents. We do so over the basic BDI architecture, but the process of deliberation and action selection is inspired by coherence maximisation inspired from Thagard’s theory of coherence.

Seen in a broader context, the theory of coherence can draw parallels with other established theories. The philosophers of science have long argued about what “claims” in a theory can be supported. Popper’s view on the progress of knowledge [9] sees falsifiability as the main driving force, and knowledge as an evolving body that follows a process in which a number of theories ‘compete’ to account for a problem situation. When a set of theories is set, falsification is then the process that makes some theories fail, while allowing others to survive. In his view survival does not mean truth but ‘fitness’ to the situation. The notion of truthlikeness is for Popper a notion of verosimilitude ($V(a) = T(a) - F(a)$) that accounts for the comparison between the truth content of theory a and the falsity content of a , which permits to rank theories. As we will see later in this paper this concept is similar to the notion of ‘strength’ of a partition in a coherence graph. Falsification of a theory can be associated to the introduction of a highly incoherent fact that will make certain statements to be removed from the accepted set of claims. Although Popper would reject a complete theory as soon as empirical evidence would go against it, Kuhn [10] would consider that scientists tolerate a certain level of anomalies (in our context a certain level of incoherence) for a long time until a revolution happens in which a complete new theory is accepted and an old one rejected. This latter phenomenon may be reproduced in our context, as we will see, by the fact that partitions in graphs can change abruptly when two theories are similarly coherent and a new experimental result is added leading to a swap in the set of accepted claims. The reconciliation point made by

Lakatos [11] would be that scientific theories contain a hard core that contains the most crucial claims of the theory plus a protective belt of auxiliary hypothesis that in case of contradiction with the facts will be modified or removed while keeping the central core, of course until a major difficulty is found that leads to a drastic change of the core. The use of degrees in claims and the algorithmic introduced in the paper will show that we might implement a similar mechanism by eliminating first the auxiliary hypothesis (those with lower degrees of belief) before removing the hard core ones (with higher probability degrees).

In the remaining of the paper, we introduce Thagard's theory of coherence, and the coherence framework used to explain action selection in Section 2. In Section 3 we explain the architecture of a coherence-driven agent and explain coherence-driven action selection. With the help of an example, we illustrate the theory in Section 4 and conclusion and future works are in Section 5.

2 Thagard's Theory of Coherence

In this section, we discuss the intuitions behind Thagard's Theory of Coherence and introduce a coherence framework based on this theory.

Paul Thagard is one of the philosophers who have attempted to introduce a computational interpretation of coherence. Thagard postulates that the theory of coherence is a cognitive theory with foundations in philosophy that approaches problems in terms of the satisfaction of multiple constraints within networks of highly interconnected elements [16; 17]. At the interpretation level, Thagard's theory of coherence is the study of associations, that is, how a piece of information influences another and how best different pieces of information can fit together. Each piece of information imposes constraints on others, the constraints being positive or negative. Positive constraints strengthen pieces of information, thereby increasing coherence, while negative constraints weaken them, thereby increasing incoherence. Hence, a coherence problem is to put together those pieces of information that have a positive constraint between them, while separating those having a negative constraint. Coherence is maximised if we obtain such a partition of information where a maximum number of constraints is satisfied.

Thagard's Formalisation

Thagard formalises coherence as follows [16]: The basic notions are that of a set of pieces of information which are represented as nodes in a graph $V = \{v_i\}$ and weighted links or constraints $E = \{\{v, w\}\}$ between these nodes. Further, some of these constraints are positive (C^+) and others negative (C^-) and associated with each constraint a number ζ which indicate the weight of the constraint. Given these, maximising coherence is formulated as the problem of partitioning V into two sets, \mathcal{A} (accepted) and \mathcal{R} (rejected), in a way that maximises compliance with the following two coherence conditions:

1. if $(v, w) \in C^+$ then $v \in \mathcal{A}$ if and only if $w \in \mathcal{A}$.
2. if $(v, w) \in C^-$, then $v \in \mathcal{A}$ if and only if $w \in \mathcal{R}$.

If $\{v, w\}$ complies with one of the above conditions, then, Thagard defines it as a satisfied constraint. Then the coherence problem is to maximise the sum of the weights of the satisfied constraints.

Thagard further proposes six main kinds of coherence: *explanatory, deductive, conceptual, analogical, perceptual, and deliberative*, each with its own array of elements and constraints. Once these elements and constraints are specified, then the algorithms that solve the general coherence problem can be used to compute coherence in ways that apply to specific domain problems.

2.1 Comparison with Other Decision Theories

Keeping Thagard's approach to coherence as maximising constraint satisfaction, we try to understand the main concept behind this theory. We associate coherence with an ever-changing system where coherence is the only property that is preserved, while everything around it changes. In cognitive terms, this would mean that, there are no beliefs nor other cognitions that are taken for granted or fixed forever. Everything can be changed and may be changed to keep coherence. We humans tend to revise or re-evaluate adherence to social norms, our plans, goals and even beliefs when we are faced with incoherence. We do not suppose that taking decisions based on coherence imply an unstable system. Our claim is based on the fact that some beliefs are more fundamental than others, in line with Lakatos. Revision of such fundamental belief is less frequent compared to other beliefs. In coherence terms, these beliefs are fundamental because they support and get support from most other cognitions and hence are in positive coherence with them. Hence, such beliefs will almost always be part of the chosen set while maximising coherence. The same is the case with other cognitions while the process of coherence maximisation further helps resolve conflicts by selecting among the best alternatives.

When applied to decision making, this means that we may not only select the set of actions to be performed to achieve certain fixed goals, but also look for the best set of goals to be pursued. Further, since coherence affects everything from beliefs to goals and actions, it may happen that beliefs contradicting a decision made are discarded. There are psychological theories such as cognitive dissonance [5] that explains this phenomenon as an attempt to justify the action chosen. Thus, with coherence we are looking at a more dynamic model of cognitions where one picks and chooses goals, actions and even beliefs to fit a grand plan of maximising coherence. In concrete terms, a highly desired state of the world (preferred in a classical sense) may get discarded in front of a less desired state of the world because it is incoherent with the rest of the beliefs, desires or intentions.

As discussed in [16], this view of decision making is very different from those of classical decision making theories where the notion of *preference* is atomic and there is no conceptual understanding of how preferences can be formed. In contrast, coherence based decision making tries to understand and evaluate these preferences from the available complex network of constraints. The assumption here is more basic because the only knowledge available to us are the various interacting constraints between pieces of information.

2.2 Coherence framework

Since we consider coherence-driven agents, in this section we summarise a generic coherence framework that will allow us to build coherence-driven agents. The framework is introduced in the work of Joseph et al [8; 7] based on Thagard's theory. It differs from other coherence-based frameworks in extending agent theories [2; 13] as in this framework coherence is treated as a fundamental property of the cognitions of an agent. Further, it is generic and fully computational. In the following we briefly introduce the necessary definitions of this framework to understand the formulation of coherence-driven action selection. The intuition behind these definitions and a few examples are given in [8; 7]. The core notion is that of a *coherence graph* whose nodes represent pieces of information and whose weighted edges represent the degree of coherence or incoherence between nodes.

Definition 2.1 A coherence graph is an edge-weighted undirected graph $g = \langle V, E, \zeta \rangle$, where

1. V is a finite set of nodes representing pieces of information.
2. $E \subseteq \{\{v, w\} | v, w \in V\}$ is a finite set of edges representing the coherence or incoherence between pieces of information, and which we shall call constraints.
3. $\zeta : E \rightarrow [-1, 1] \setminus 0$ is an edge-weighted function that assigns a negative or positive value to the coherence between pieces of information, and which we shall call coherence function.

Every coherence graph is associated with a number called the *coherence of the graph*. Based on Thagard's formalism, this can be calculated by partitioning the set of nodes V of the graph in two sets, \mathcal{A} and $V \setminus \mathcal{A}$, where \mathcal{A} contains the accepted elements of V , and $V \setminus \mathcal{A}$ contains the rejected ones. The aim is to partition V such that a maximum number of constraints is satisfied, taking their values into account. A constraint is satisfied only if it is positive and both the end nodes are in the same set, or negative and the end nodes are in complementary sets. The following definitions help clarify this idea.

Definition 2.2 Given a coherence graph $g = \langle V, E, \zeta \rangle$, and a partition $(\mathcal{A}, V \setminus \mathcal{A})$ of V , the set of satisfied constraints $C_{\mathcal{A}} \subseteq E$ is given by

$$C_{\mathcal{A}} = \left\{ \{v, w\} \in E \mid \begin{array}{l} v \in \mathcal{A} \text{ iff } w \in \mathcal{A} \text{ when } \zeta(\{v, w\}) > 0 \\ v \in \mathcal{A} \text{ iff } w \notin \mathcal{A} \text{ when } \zeta(\{v, w\}) < 0 \end{array} \right\}$$

All other constraints (in $E \setminus C_{\mathcal{A}}$) are said to be unsatisfied.

Definition 2.3 Given a coherence graph $g = \langle V, E, \zeta \rangle$, the strength of a partition $(\mathcal{A}, V \setminus \mathcal{A})$ of V is given by

$$\sigma(g, \mathcal{A}) = \frac{\sum_{\{v, w\} \in C_{\mathcal{A}}} |\zeta(\{v, w\})|}{|E|}$$

Notice that, by Definitions 2.2 and 2.3,

$$\sigma(g, \mathcal{A}) = \sigma(g, V \setminus \mathcal{A}) \quad (1)$$

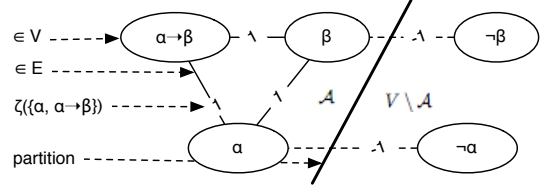


Figure 1: A typical coherence graph with a coherence maximising partition

Definition 2.4 Given a coherence graph $g = \langle V, E, \zeta \rangle$ and given the strength $\sigma(g, \mathcal{A})$, for all subsets \mathcal{A} of V , the coherence of g is given by

$$\kappa(g) = \max_{\mathcal{A} \subseteq V} \sigma(g, \mathcal{A})$$

If for some partition $(\mathcal{A}, V \setminus \mathcal{A})$ of V , the strength of the partition is maximal (i.e., $\kappa(g) = \sigma(g, \mathcal{A})$) then the set \mathcal{A} is called the accepted set and $V \setminus \mathcal{A}$ the rejected set of the partition. A typical coherence graph is as shown in Figure 1.

Due to Equation 1, the accepted set \mathcal{A} is never unique for a coherence graph. Moreover, there could be other partitions that generate the same value for $\kappa(g)$. Here we mention a few criterias to select an accepted set among the alternatives. If $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$ are sets from all those partitions that maximise coherence of the graph g , based on one of Thagard's principles (which we will formalise in the next definition) on deductive coherence[16] that *intuitively obvious propositions have an acceptability on their own*, we say an accepted set is the one in which the intuitively obvious propositions belong. Further, the coherence of the sub-graphs $(g|_{\mathcal{A}_i}, i \in [1, n])$ gives us an indication of how strongly connected they are. The higher the coherence, the more preferred the corresponding accepted set. And lastly, an accepted set with more number of elements should be preferred to another with less.

We now need a way in which the coherence graphs just defined can be constructed. That is, we need to define function ζ . As the nature of relationship between two pieces of information (corresponding to the different types of coherence as mentioned in the introduction) can vary greatly, we do not have one unique coherence function. That is, in an explanatory coherence, two pieces of information are coherent when they are related by an explanation. Thagard proposes certain principles to characterise coherence in each of the different types. Here we define one such coherence function which is inspired from Thagard's principles of *deductive coherence*.

Thagard's principle mainly states that *a proposition coheres with propositions that are deducible from it, propositions that are used together to deduce something cohere with each other, the more hypotheses it takes to deduce something, the less the degree of coherence, contradictory propositions are incoherent with each other*¹. Since some of these principles make sense only in the context of a theory presentation,

¹here we do not formalise the principle that *intuitively obvious propositions have a degree of acceptability on their own*. This we keep it as a disambiguation criteria to select among accepted sets.

we assume a theory presentation \mathcal{T} in a multi-valued propositional logic while formalising these principles. We use a multi-valued logic to model uncertainty in agents, though Thagard's principles, we assume, are based on a boolean world. We formalise Thagard's principles in terms of a *support function* $\eta_{\mathcal{T}}$ which extract a coherence value between two nodes if either one implies the other, or together they are used to imply a third node (assuming some sort of deduction theorem such as $\mathcal{T}, \alpha \vdash \beta$ implies $\mathcal{T} \vdash \alpha \rightarrow \beta$). We also normalise the values between $[-1, 1]$.

Definition 2.5 *Let L be the set of all propositional sentences of a multi-valued propositional logic. Let $\mathcal{T} \subseteq L$ be a finite theory presentation and $\Gamma \subseteq \mathcal{T}$ and $\gamma \in \bar{L}$. A support function $\eta_{\mathcal{T}} : L \times L \rightarrow [-1, 1]$ with respect to \mathcal{T} is given by*

$$\eta_{\mathcal{T}}(\{\alpha, \beta\}) = \begin{cases} \max \left\{ \begin{array}{l} \max \left\{ \frac{2 \cdot F_{\rightarrow}(\rho(\alpha), \rho(\beta)) - 1}{|\Gamma|} \mid \right. \\ \left. \exists \Gamma \subseteq \mathcal{T} : \Gamma, \alpha \vdash \beta ; \alpha \not\vdash \beta \right\}, \\ \max \left\{ \frac{2 \cdot F_{\rightarrow}(\rho(\alpha), F_{\rightarrow}(\rho(\beta), \rho(\gamma))) - 1}{|\Gamma| + 1} \mid \right. \\ \left. \exists \Gamma \subseteq \mathcal{T} : \Gamma, \alpha, \beta \vdash \gamma ; \alpha, \beta \not\vdash \gamma \right\} \end{array} \right\} \\ \text{undefined} & \text{otherwise} \end{cases}$$

where F_{\rightarrow} is the truth connective defined for L and $\rho(\alpha)$ gives the truth value of α .

Thagard in his principles emphasises the fact that, though a coherence value can be derived from the underlying implication relation, coherence functions are always symmetric. Due to this, even if there may only be a deductive relation in one direction, there will be a deductive coherence in both directions. Hence, we define the deductive coherence between two propositions as the value of the stronger $\eta_{\mathcal{T}}$ values.

Definition 2.6 *Let L be the set of all propositional sentences of a multi-valued propositional logic. Let $\mathcal{T} \subseteq L$ be a finite theory presentation and let $\eta_{\mathcal{T}} : L \times L \rightarrow [-1, 1]$ be a support function. A deductive coherence function $\zeta_{\mathcal{T}} : L \times L \rightarrow [-1, 1] \setminus \{0\}$ with respect to \mathcal{T} is a partial function given by: For any pair (α, β) of formulas in L ,*

$$\zeta_{\mathcal{T}}(\{\alpha, \beta\}) = \begin{cases} \max(\eta_{\mathcal{T}}(\alpha, \beta), \eta_{\mathcal{T}}(\beta, \alpha)) \\ \text{if } \eta_{\mathcal{T}}(\alpha, \beta) \text{ and } \eta_{\mathcal{T}}(\beta, \alpha) \text{ defined, } \neq 0 \\ \eta_{\mathcal{T}}(\alpha, \beta) \\ \text{if } \eta_{\mathcal{T}}(\alpha, \beta) \text{ defined and } \neq 0 \\ \text{and } \eta_{\mathcal{T}}(\beta, \alpha) \text{ undefined or } = 0 \\ \eta_{\mathcal{T}}(\beta, \alpha) \\ \text{if } \eta_{\mathcal{T}}(\beta, \alpha) \text{ defined and } \neq 0 \\ \text{and } \eta_{\mathcal{T}}(\alpha, \beta) \text{ undefined or } = 0 \\ \text{undefined} \\ \text{if } \eta_{\mathcal{T}}(\alpha, \beta) \text{ and } \eta_{\mathcal{T}}(\beta, \alpha) \text{ are undefined} \\ \text{or } = 0 \end{cases}$$

Note that both the support function and the deductive coherence function are partial functions. This is because we interpret zero coherence as the propositions not being related.

3 Coherence-driven Agent Architecture

A *coherence-driven agent* is an agent which always takes an action based on maximisation of coherence of its cognitions,

norms and other social commitments. Further, these are cognitive agents based on BDI theory [14] and are modeled as a multi-context architecture (developed by Casali et al. [3]), which consists of a set of contexts and a set of bridge rules between contexts. Each context has its own language, logic and theory expressed as coherence graphs. Bridge rules turn formulae derivable in one or more contexts into premises for derivations for another context. We assume that each agent has its beliefs, desires, and intentions stored in its belief context C_B , desire context C_D , and intention context C_I .

3.1 Cognitive Contexts

Here we briefly describe how a belief context C_B is defined while desire C_D and intention C_I contexts are similar [8; 3]. C_B consists of a belief logic and a theory \mathcal{T}_B of the logic expressed as a coherence graph.

A belief logic \mathcal{K}_B consists of a belief language, a set of axioms and a deductive relation defined on the belief logic $\langle L_B, A_B, \vdash_B \rangle$. The belief language L_B is defined by extending the classical propositional language L defined upon a countable set of propositional variables PV and connectives (\neg, \rightarrow) . L is extended with a fuzzy unary modal operator B . The modal language L_B is built from the elementary modal formulae $B\varphi$ where φ is propositional, and truth constants r , for each rational $r \in \mathbb{Q} \cap [0, 1]$, using the connectives of Łukasiewicz many-valued logic. If φ is a proposition in L , the intended meaning of $B\varphi$ is that “ φ is believable”. A modal many-valued logic based on Łukasiewicz logic is used to formalise \mathcal{K}_B ².

Definition 3.1 [3] *Given a propositional language L , a belief language L_B is given by:*

- If $\varphi \in L$ then $B\varphi \in L_B$
- If $r \in \mathbb{Q} \cap [0, 1]$ then $\bar{r} \in L_B$
- If $\Phi, \Psi \in L_B$ then $\Phi \rightarrow_L \Psi \in L_B$ and $\Phi \& \Psi \in L_B$ (where $\&$ and \rightarrow_L correspond to the conjunction and implication of Łukasiewicz logic)

We call \mathcal{T}_B a theory in the language L_B .

Other Łukasiewicz logic connectives for the modal formulae can be defined from $\&$, \rightarrow_L and $\bar{0}$: $\neg_L \Phi$ (defined as $\Phi \rightarrow_L \bar{0}$). Formulae of the type $\bar{r} \rightarrow_L \Psi$ (the probability of φ is at least r) will be denoted as (Ψ, r) .

The axioms A_B of \mathcal{K}_B are:

1. All axioms of propositional logic.
2. Axioms of Łukasiewicz logic for modal formulas (for instance, axioms of Hájek's Basic Logic (BL) [6] plus the axiom: $\neg\neg\Phi \rightarrow \Phi$.)
3. Probabilistic axioms, given $\varphi, \psi \in L$:
 - $B(\varphi \rightarrow \psi) \rightarrow_L (B\varphi \rightarrow B\psi)$
 - $B\varphi \equiv \neg_L B(\varphi \wedge \neg\psi) \rightarrow_L B(\varphi \wedge \psi)$

The deduction rules defining \vdash_B of \mathcal{K}_B are Modus ponens and Necessitation for B (from φ derive $B\varphi$).

Note that the truth function $\rho : L_B \rightarrow [0, 1]$ is defined by means of the truth-functions of Łukasiewicz logic and the probabilistic interpretation of beliefs as follows:

²We could use other logics as well by replacing the axioms.

- $\rho((B\varphi, r))^3 = r$ for all $r \in \mathbb{Q} \cap [0, 1]$
- $\rho(\varphi \& \psi) = \max(\rho(\varphi) + \rho(\psi) - 1, 0)$ for all $\varphi, \psi \in L_B$
- $\rho(\varphi \rightarrow_L \psi) = \min(1 - \rho(\varphi) + \rho(\psi), 1)$ for all $\varphi, \psi \in L_B$

Then a coherence graph over beliefs is defined over the belief logic \mathcal{K}_B as follows:

Definition 3.2 Given a belief logic $\mathcal{K}_B = \langle L_B, A_B, \vdash_B \rangle$ where L_B is a belief language, A_B are a set of axioms and \vdash_B are a set of deduction rules, a belief coherence graph $g_B = \langle V_B, E_B, \zeta_B \rangle$ is a coherence graph defined over \vdash_B and a finite theory \mathcal{T}_B of L_B such that:

- $V_B \subseteq \mathcal{T}_B$
- E is a set of subsets of 2 elements of V_B
- $\zeta_{\mathcal{T}_B}$ is defined over \vdash_B and \mathcal{T}_B .

A belief coherence graph exclusively represents the graded beliefs of an agent and the associations among them. A desire coherence graph (g_D), and an intention coherence graph (g_I) over logics L_D , and L_I are similar.

3.2 Bridge Rules

Bridge rules are inference rules of the form $b = \frac{C_1:\psi, C_2:\psi}{C_3:\psi}$ whose premises and conclusion are labelled formulas where the labels denote the contexts they are taken from. They carry inferences between theories of different logics. Since our theories become coherence graphs, we need two functions to emulate the execution of bridge rules over coherence graphs. If \mathcal{G} denote the set of all coherence graphs, then a graph node extension function ($\varepsilon : \mathcal{G}^n \rightarrow \mathcal{G}^n$) takes into account the influence of graphs (theories) on each other. An edge extension function ($\iota : \mathcal{G}^n \rightarrow \mathcal{G}$) joins a set of graphs by adding edges between the nodes participating in the inference. Since we treat bridge rules similar to any other implication relations, we use Definition 2.6 itself to calculate the coherence values on these edges. We now illustrate the concept of bridge rules when the contexts are coherence graphs (the formal definitions can be found in [8]).

Example 3.1 Let's assume, for instance, that an agent wants it to be the case that whenever it has an intention ($I\varphi, r$) in the intention graph (a formula in the theory \mathcal{T}_I), then the corresponding belief ($B\varphi, r$) is inferred in the belief graph (added to the theory \mathcal{T}_B). i.e., Given a bridge rule $b = \frac{C_B:(B\psi, r), C_D:(D\psi, s)}{C_I:(I\psi, \min(r, s))}$ where contexts C_B, C_D , and C_I have the coherence graphs g_B, g_D and g_I associated with them respectively and given $(B\psi, 0.95) \in g_B$, $(D\psi, 0.95) \in g_D$ function ε adds a node ($I\psi, 0.95$) to g_I .

Let's further assume that our agent further wants it to be the case that, the belief and the intention nodes are related and have a positive coherence between them. The edge extension function ι joins the graphs g_B, g_D and g_I associated with the contexts in the bridge rule by adding the edges $\{(I\psi, 0.95), (B\psi, 0.95)\}, \{(I\psi, 0.95), (D\psi, 0.95)\}$ with coherence values equal to $\frac{2 \cdot \rho((I\psi, 0.95)) - 1}{1} = 0.9$ from Definition 2.5.

³ $(B\varphi, r) \equiv \bar{r} \rightarrow B\varphi$

3.3 Architecture

Figure 2 shows the architecture of a coherence-driven agent. At any time, it can either perceive the environment (updates beliefs) or make a decision about a future action. In the event of a new information, an agent re-evaluates its theory, hence recomputes both the coherence graphs and the coherence maximising partition. If the new information falls in the accepted set then it reinforces the theory and the theory becomes more coherent. However, if it falls in the rejected set, then it contradicts some of the elements of the accepted theory to make the theory more coherent again in line with Lakatos. In the process, some of the existing elements may move from accepted to rejected or vice versa. An agent always bases its decisions on the accepted theory.

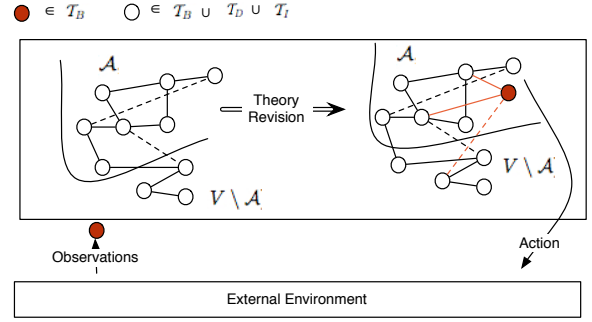


Figure 2: A coherence-driven agent architecture

3.4 Coherence-driven Action Selection

As discussed in the introduction, the philosophy behind a coherence-driven action selection is substantially different from other typical goal-driven approaches to action selection. A BDI-based agent, at anytime selects a goal to pursue, and looks for what actions would satisfy that goal. It is argued that, this would reduce the attention problem of the agent, giving it a stable behaviour. However, as argued in the introduction, this has many difficulties such as incorporating new perceptions, analysing conflicts between goals, and analysing feasibility of actions to achieve goals. This is due to the philosophical grounding of the theory for which the basis for an action is the expectation of a desired outcome.

As suggested in the introduction, coherence-driven reasoning offers a more holistic view on action selection. The philosophical theories of action suggest that an agent is influenced by a reason, and his action is consequently performed for that reason, when he is influenced by a representation of the action that makes it intelligible to him. Naturally, this representation may make the action intelligible precisely by setting it in the context of his desires and expectations, but his reason for action consists in this cognitively attractive representation of it rather than in the desires and expectations to which it alludes. A reason is a rationale, in the light of which an action makes sense to an agent, and promoting a desired outcome is one such rationale [18]. The coherence-driven approach we propose here attempts to capture this representation, which gives the agent the necessary rationale for action.

At any time a coherence-driven agent selects the most preferred action from its current accepted set of a coherence maximising partition. Any external stimuli interrupts the action selection process and forces the representation of the cognitions to go through a re-evaluation of coherence, resulting in some of the currently accepted cognitions to be rejected and vice versa. A procedure that a typical coherence-driven agent follows is outlined in the following.

Given the current coherence graphs g_B , g_D , and g_I and their composition ς_g and an external stimuli $(K\varphi, r)$ where $K \in \{B, D, I\}$

```

1: if  $(K\varphi, r)$  then
2:    $v := (K\varphi, r)$ 
3:    $V_K := V_K \cup \{v\}$ 
4:   for all  $w \in V_K$  do
5:     compute  $\zeta(\{v, w\})$  using Definition 2.6.
6:     if  $\zeta(\{v, w\})$  is defined then
7:        $E_K := E_K \cup \{\{v, w\}\}$ 
8:     end if
9:   end for
10:  compute a composite coherence graph  $\varsigma_g$  as in [8] and Example 3.1.
11:  for all  $(\mathcal{A}_i, V \setminus \mathcal{A}_i), \mathcal{A}_i \subseteq V$  do
12:    calculate  $\sigma(\varsigma_g, \mathcal{A}_i)$  using Equation 2.3
13:  end for
14:   $\kappa := \kappa(\varsigma_g)$  using Equation 2.4
15:   $\mathcal{A} := \mathcal{A}_i | \max(\sigma(\varsigma_g, \mathcal{A}_i))$ 
16: end if
17:  $current\_action := \max_r \{(I\varphi, r) | (I\varphi, r) \in \mathcal{A}\}$ 

```

Line 1 checks for external stimuli. If there are any, then, lines from 2 to 9 updates the graphs by incorporating the stimuli and its influences on existing elements of the observed cognition. Line 10 builds up the reasoning across contexts by composing the coherence graphs. Lines from 11 to 14 determines the coherence maximising partition. This is done by first computing the strength of each partition using the function σ and choosing the partition $(\mathcal{A}, V \setminus \mathcal{A})$ for which $\sigma(g, \mathcal{A})$ is maximal. This part of the algorithm only gives the simplest solution, however, finding a maximising partition of a weighted graph is known to be an NP-complete problem. There are approximation algorithms exist to find the solution to this problem such as max-cut, neural network based algorithms. Line 17 determines the current action by selecting the action from \mathcal{A} which has the highest preference.

4 Example

We consider a simple example to show how action selection works in our architecture. A coherence-driven robotic agent wants to choose between a set of possible actions(intentions) corresponding to a set of desires (goals) it has. The scenario is modelled like a grid in which at each cell the robot can chose between two possible actions: “plug” to restore its energy or “move” to earn points. It is further assumed that at every cell in the grid, it is possible to perform both actions. With every move the robot gains a point. Finally, the robot is equipped with an energy sensor, which measures the remaining energy at every time point, which influences the choices of the robot. The results are based on an implementation of a heuristic-

based polynomial-time approximation algorithm to compute partitions and their corresponding coherences.

Since coherence maximisation dynamically choses the most coherent partition, the robot at any instant choses the set of beliefs, desires and intentions (actions) that it wants to pursue. There is one single persistent desire for the robot, which is to earn points. It has certain domain knowledge which indicates how to get its desire satisfied. This domain knowledge is encoded as a belief $(B(move \rightarrow points), 1)$ says that a move will fetch a point with a confidence degree 1. A bridge rule b_1 is used to reason with the beliefs and desires.

$$b_1 = \frac{C_B : (B(p \rightarrow q), \alpha), C_D : (Dq, \beta)}{C_D : (Dp, \min(\alpha, \beta))}$$

b_1 injects a new desire p given the desire of q and a belief that p facilitates q with appropriate degrees. Further, using b_1 and the belief that *having energy enables move*, i.e., $(B(energy \rightarrow move), 1)$, a new desire to have “energy” is generated. A third desire to “plug” is generated using the bridge rule and the belief that *plugging gives energy*, i.e., $(B(plug \rightarrow energy), 1)$. The chain of desires and their coherence links are illustrated in Figure 3.

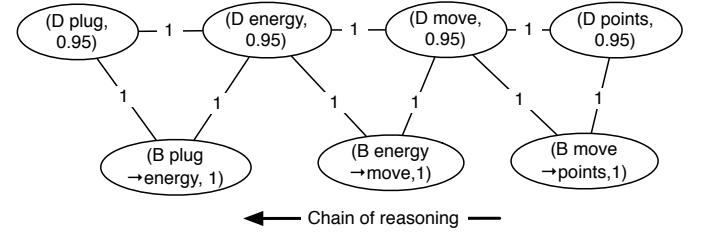


Figure 3: How one desire triggers another

The robot uses a second bridge rule b_2 that states that every desire with a corresponding belief that the desire is achievable, generates a corresponding intention (realistic agent).

$$b_2 = \frac{C_B : (Bp, \alpha), C_D : (Dp, \beta)}{C_I : (Ip, \min(\alpha, \beta))}$$

Using this rule, it has the intention to move, intention to have energy and intention to plug. Further, as in the case of desires, the intention to move is connected to the intention to have energy using another bridge rule b_3 which is used to reason across beliefs and intentions.

$$b_3 = \frac{C_B : (B(p \rightarrow q), \alpha), C_I : (Dq, \beta)}{C_I : (Dp, \min(\alpha, \beta))}$$

Note that bridge rules b_1 and b_3 are very similar and motivated from the well known practical syllogism, “If I want q and p realises q , then I should intend to do p ”. Using b_3 , we have that *the belief* $(B(energy \rightarrow move), 1)$, $(Ienergy, x)$ and $(Imove, x)$ are coherently related. The same is true of $(Ienergy, x)$ and $(Iplug, x)$. Hence, similar to desires, a chain of intentions and their coherence links are generated (Figure 4).

As the only sensor for the robot (other essential sensors ignored) relevant to the problem is the *energy_sensor*

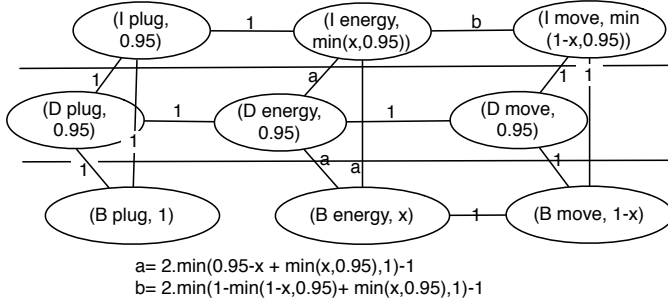


Figure 4: Desires and intentions trigger other intentions

(e_s), at every time point a few of the cognitions gets affected due to the changes in sensor readings. That is, we take that the grade on the belief that the quantity of energy needed changes inversely to the value of e_s . Further, the belief that move is possible changes proportionally with the value of e_s , using the modes ponens as $(B(\text{energy} \rightarrow \text{move}), 1)$, $(B(\text{energy}, x) \rightarrow (B(\text{move}, 1 - x))$. Finally, as it is assumed that the robot can perform only one action at a time, the essential conflict between intentions to “move” and “plug” are expressed as $(I(\text{move}, x) \Leftrightarrow (I(\neg \text{plug}, x))$.

4.1 Action Selection

Given our robotic agent as described, we now pose the problem of action selection. That is, the robot has to decide what action to perform at every time point. We say an “energy-cycle” is the time between two consecutive “plug” actions. We take different energy levels and determine both the coherence of the robot and the coherence-driven choice of action. To understand how the coherence graph would look like, we show the graph with the partition when the $x = 0$ (just plugged) in Figure 5 and when $x = 1$ (there is no energy left) in Figure 6. In the case $x = 0$, there is a clear partition with the only intention selected is $(I(\text{move}, 0.95))$. Hence it is absolutely certain that, the robot should chose to move and earn points. The coherence of the graph for this partition is 0.6533.

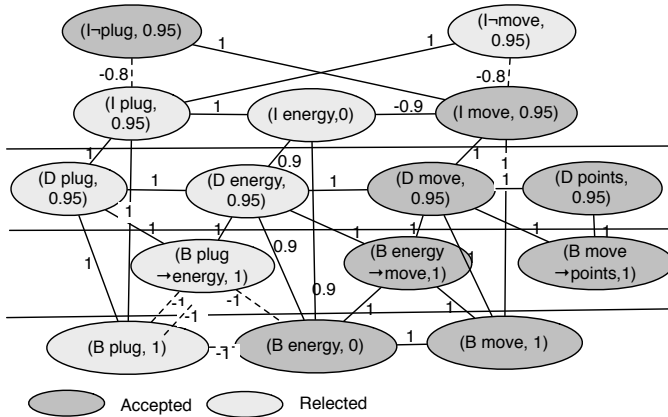


Figure 5: $x=0$ (Robot has maximum energy)

In the case $x = 1$ (Figure 6), there are no incoherence relations in the graph. However, it is due to the fact that, we are deriving the belief about move from the belief about energy needs. Though every node is part of the accepted set, notice that the node with the highest grade will be pursued. In this case the choice to plug will be pursued. The coherence of the graph for this partition is 0.97777. The increase in coherence is due to the fact that, there are no incoherence experienced by the robot and hence all nodes are accepted. Most of the individual coherence values are equal to the maximum possible ($= 1$).

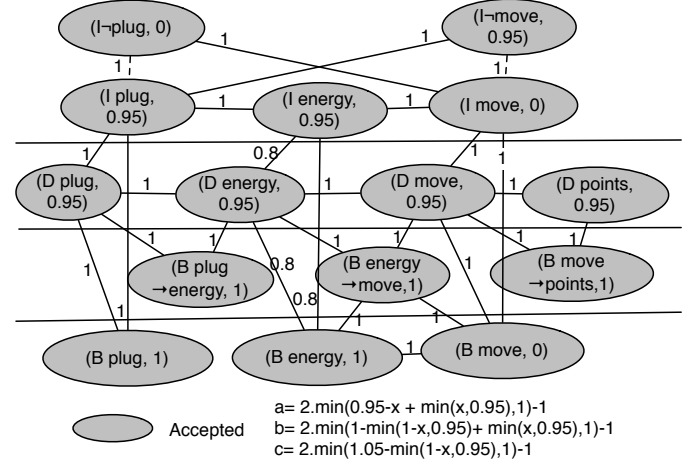


Figure 6: $x=1$ (Robot has no energy)

These are the two extreme cases, and now we plot the behaviour of the robot in terms of the choice of action, and the variation in coherence values at different energy levels in two energy-cycles. As seen in the coherence graphs in Figures 5 and 6, when the energy requirement is 0 or close to 0, the robot has some incoherences and selects only few of the cognitions as accepted. However, as the energy requirement increase, these incoherences disappear (due to the decreasing intention for action “move”) and hence the robot becomes increasingly coherent with the action to plug.

Another graph which is interesting is the energy levels versus the choice of action as in Figure 8. This shows the expected actions of the robot at different energy levels. When the energy need is in the range $[0, 0.5]$ the robot choses to move. However, if the energy need is in the range $[0.6, 1]$, then the robot choses to plug and restore the energy. Then, its clear that as soon as the energy need raises to 0.6, the robot take the action to plug. Thus, the energy need never raises to a point beyond 0.6 (conservative behavior). Hence, the actual behavior of the robot will be a repeating sequence of $\{\text{Move, Move, Move, } \dots, \text{Plug, Move, Move, } \dots\}$, as the intuition would make us expect.

5 Discussion and Futurework

In this paper, we have introduced an alternative approach to action selection based on coherence maximisation. The interesting aspects of this approach over more traditional BDI

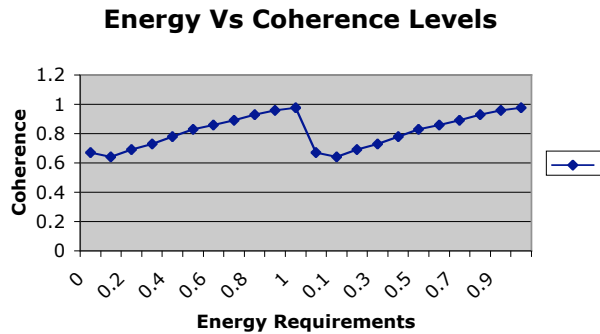


Figure 7: Variation in coherence with different energy levels

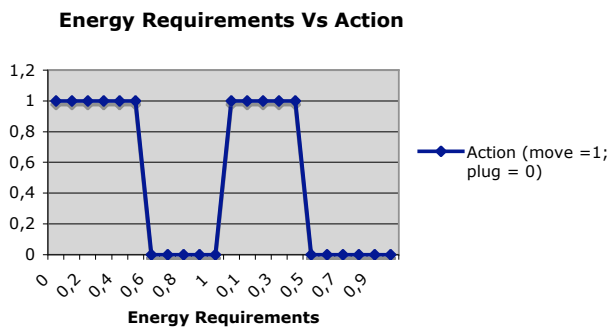


Figure 8: Action selection with different energy levels

approaches are that it takes a dynamic view of agent cognitions, can detect and resolve conflicts among cognitions, can perform uncertainty reasoning and can reason at a global level while also fully integrated into the BDI representation. Since we have discussed related work in the course of presenting the paper, we here make a brief comment on one related work, which is the only work known to us that uses coherence for agent reasoning. While the work of Pasquier et. al [13; 12], introduced coherence based reasoning in agents, there are significant differences with our proposal. In their work, coherence is like a utility maximising function, which is used to prioritise the intentions (dialogue moves), whereas reasoning about beliefs and desires are using the traditional BDI approach. This we imagine will retain all the difficulties we mentioned in the introduction. Another important difference is that, while we show how coherence can be computed using Thagard's principles, such mechanisms are missing from their approach.

In the future work, we plan to incorporate the representation of plans and study how plans can be included in the coherence maximising process. Further, we plan to explore the possibilities of evaluating our framework using empirical and mathematical proof.

References

- [1] M. E. Bratman. *Intention, Plans, and Practical Reason*. CSLI publications, 1987.
- [2] J. Broersen, M. Dastani, J. Hulstijn, Z. Huang, and L. van der Torre. The BOID architecture: conflicts between beliefs, obligations, intentions and desires. In *AGENTS '01*. ACM, 2001.
- [3] A. Casali, L. Godo, and C. Sierra. Graded BDI models for agent architectures. In *Lecture Notes in Computer Science*, volume 3487, 2005.
- [4] M. Dastani, F. de Boer, F. Dignum, and J.-J. Meyer. Programming agent deliberation: an approach illustrated using the 3apl language. In *AAMAS '03*, pages 97–104. ACM, 2003.
- [5] L. Festinger. *A theory of cognitive dissonance*. Stanford University Press, 1957.
- [6] P. Hájek. Metamathematics of fuzzy logic. In *Trends in Logic*, volume 4, 1998.
- [7] S. Joseph, C. Sierra, and M. Schorlemmer. A coherence based framework for institutional agents. In *Lecture Notes in Computer Science*, volume 4870, 2007.
- [8] S. Joseph, C. Sierra, M. Schorlemmer, and P. Dellunde. Formalising deductive coherence: An application to norm evaluation. In *Normas'08(Extended version)*, 2008.
- [9] K. Popper. *Conjectures and Refutations: The Growth of Scientific Knowledge*. Basic Books, 1962.
- [10] T. S. Kuhn. *Conjectures and Refutations: The Growth of Scientific Knowledge*. Chicago: University of Chicago Press, 1962.
- [11] I. Lakatos. *Conjectures and Refutations: The Growth of Scientific Knowledge*. Cambridge University Press, 1976.
- [12] P. Pasquier, N. Andrillon, M.-A. Labrie, and B. Chaib-draa. An exploration in using cognitive coherence theory to automate bdi agents communicational behavior. In *Advances in Agent Communication*. Springer, 2004.
- [13] P. Pasquier and B. Chaib-draa. The cognitive coherence approach for agent communication pragmatics. In *Second International Joint Conference on Autonomous Agents and Multiagent Systems*, 2003.
- [14] A. S. Rao and M. P. Georgeff. Bdi agents: From theory to practice. In *ICMAS-95, First International Conference on Multi-Agent Systems: Proceedings*, pages 312–319. MIT Press, 1995.
- [15] Y. Shoham. *Agento: a simple agent language and its interpreter*. In *Proceedings of AAAI*, 1993.
- [16] P. Thagard. *Coherence in Thought and Action*. MIT Press, 2002.
- [17] P. Thagard. *Hot Thought*. MIT Press, 2006.
- [18] J. D. Velleman. *Self to Self: Selected Essays*. Cambridge University Press, 2005.