



King's Research Portal

DOI:

[10.1007/978-3-319-59294-7_21](https://doi.org/10.1007/978-3-319-59294-7_21)

Document Version

Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Santacà, K., Cristani, M., Rocchetto, M., & Viganò, L. (2017). A Topological Categorization of Agents for the Definition of Attack States in Multi-Agent Systems. In *Proceedings of EUMAS 2016: Part of the Lecture Notes in Computer Science book series* (Vol. 10207, pp. 261-276) https://doi.org/10.1007/978-3-319-59294-7_21

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

A Topological Categorization of Agents for the Definition of Attack States in Multi-Agent Systems

Katia Santacà¹, Matteo Cristani¹, Marco Rocchetto², and Luca Viganò³

¹ Dipartimento di Informatica, Università di Verona, Italy

² iTrust, Singapore University of Technology and Design

³ Department of Informatics, King's College London, UK

Abstract. We propose a topological categorization of agents that makes use of the multiple-channel logic (MCL) framework, a recently developed model of reasoning about agents. We firstly introduce a complete formalization of prejudices on agents' attitudes and propose an extension of the rules of the MCL framework. We then use RCC-5 (the Region Connection Calculus) to categorize different agents in Multi-Agent Systems (MAS) based on the collaboration, competence, and honesty of agents. We discuss the possibility of using RCC-3 and RCC-8 and generalize our results to define an upper bound in the number of different types of agents in MAS. Finally, we apply our topological categorization to a specific MAS that describes a Cyber-Physical System, for which we define, categorize and discuss the resulting attack states.

1 Introduction

Much effort has been devoted to the characterization of different agents in *Multi-Agent Systems (MAS)*, ranging from works that employ Dynamic Epistemic Logic and Public Announcement Logics (PAL) [13] to more recent approaches such as [1]. These works have studied an agent's beliefs and announcements, typically under the assumption that agents are always truthful and sincere. However, as discussed in, e.g., [2], this assumption is an oversimplification since most MAS contain a number of agents that are clearly neither sincere nor truthful. This is, for instance, the case in the systems that are typically considered in the security research community, where dishonest (and thus neither sincere nor truthful) agents are used to formalize attacks to the systems under consideration. As a result, a number of research paper have focused their attention to spotting unintended or even malicious behavior in MAS. We are specifically focused on Cyber-Physical Systems (CPS), as samples of such problems, as reported widely in [11,6], and specifically for an agent-based model of CPS in [7].

Distinguishing between the different types of agents in a MAS is a difficult task. This is witnessed by the fact that although a characterization of agents would obviously play a crucial role in the understanding of different aspects and facets in MAS, a proper definition is still missing. General problems with agency

and norms, namely social regulations, are presented in [3], where many problems are discussed and left open. Further investigations, including those in Public Announcement Logic [1] have devised a pathway to follow, with many problems in the definitions yet open. A step in this direction has been done in [2], which introduced a general logical framework called *Multiple Channel Logic Framework (MCL)*. However, the focus of [2] is on the definition of the framework and little attention is paid to the definition of a general categorization of all the possible agents that could be defined using MCL.

The overall goal of this paper is the definition of a general categorization of agents, based on MCL. We focus, in particular, on the application of MAS for reasoning about security systems, such as *Cyber-Physical Systems*. More specifically, our contribution is three-fold:

1. We define a topological categorization of agents in MAS, obtaining 50 new rules in the MCL framework.
2. We identify a theoretical limit to the maximum number of different types of agents in a MAS (defined using MCL).
3. As an example of a concrete application, we apply our topological categorization to define attack states for a MAS that describes a general CPS. Our case study ultimately allows us to show that our categorization of agents can be used to reason on the security of CPS and, more generally, MAS.

We proceed as follows. In Section 2, we provide background on the MCL framework that we have employed as a basis of our categorization of agents. In Section 3, we define agents and summarize the Region Connection Calculus. In Section 4, we propose a categorization of agents for MAS by extending the MCL framework. We then define an upper bound on the number of different types of agents in Section 5, and in Section 6, we apply our categorization to the security of CPS. In Section 7, we take some conclusions.

2 Background: The MCL framework

In this section, we summarize the main features of the Multiple Channel Logic Framework MCL of [2], which provides the basis for our work in this paper.

2.1 Announcements, Beliefs and Facts

MCL is a logical framework that is able to relate announcements, agents' beliefs, and true statements on multiple communication channels, where the channels of MCL are logical spaces in which agents make *public* announcements (private channels are out of the scope of MCL). More specifically, MCL is a three-layered, labeled, modal logic framework:

- The first layer is a propositional calculus that is used to express what agents share, i.e., the logical representation of an assertion.

- The second layer is a multi-modal calculus with three different modalities: B (*belief*), which allows one to assert that an agent believes in a proposition, and T_{\square} and T_{\diamond} to state that a given proposition is *asserted* by an agent respectively in every channel or at least one channel.
- The third layer is for *agent tagging*, which defines prejudice about communicative attitudes of agents (see Section 2.3 for more details).

Propositional formulae in the first layer of MCL are of the form

$$\varphi := A \mid \neg\varphi \mid \varphi \wedge \varphi \mid \varphi \vee \varphi,$$

where A denotes a propositional letter, and \neg , \wedge and \vee are the standard connectives for negation, conjunction and disjunction, respectively.

Modal formulae in the second layer of MCL are of the form

$$\mu := B[\lambda : \varphi] \mid T_{\square}[\lambda : \varphi] \mid T_{\diamond}[\lambda : \varphi] \mid \sim\mu$$

where φ denotes a propositional formula, λ an agent and \sim a negation.

$B[\lambda : \varphi]$ intuitively means that the agent λ believes in the formula φ . Note that, as is standard, an agent might believe a false formula.

$T_{\square}[\lambda : \varphi]$ intuitively means that the agent λ announces φ in every channel. More formally, when λ announces φ in a channel C , then he announces φ in any channel C' that is accessible from C . T_{\diamond} denotes that the agent λ announces φ at least in one channel. In fact, the semantics of MCL relates the notion of accessibility to the notion of *observation*: a channel C' is accessible from a channel C when the observer of C also observes C' . We won't however go into the details of the Kripke-style semantics of MCL, which is given in [2], where the soundness and completeness of MCL are proved.

We can then define the following three sets with respect to an agent λ :

- *Announcements* $\mathbb{A}_{\lambda} = \{\varphi.T_{\diamond}[\lambda : \varphi]\}$ is the set of formulae announced by λ in one or more channels.
- *Beliefs* $\mathbb{B}_{\lambda} = \{\varphi.B[\lambda : \varphi]\}$ is the set of the formulae believed to be true by λ .
- *Facts* $\mathbb{F} = \{\varphi \mid \sigma(\varphi) = \top\}$ is the set of *true* formulae.

2.2 Assumptions on the Agents in MCL

MCL imposes the following assumptions on the agents:

- *Atemporal channels*: announcements are made in a channel and hold forever.
- *Belief revision*: if an agent makes two opposite announcements in the same channel, then the agent has changed his point of view.
- *Coherent agents*: an agent makes coherent announcements in a single channel, although he might make opposite announcements in a different channel.
- *Consistent agents*: an agent either believes in the truthfulness of a statement or in the truthfulness of the opposite statement, but not in both at once.
- *No beliefs*: if an agent does not assert something, this doesn't imply that the agent believes the opposite.
- *Provable facts*: there exist provable facts that are not matter of opinions.

Table 1. Description of different agent types in MCL.

Agent type	Notation	Description
<i>Weakly Collaborative</i>	$+(W_{CI})\lambda$	λ announces his beliefs in at least one channel
<i>Strongly Collaborative</i>	$+(S_{CI})\lambda$	λ announces his beliefs in every channel
<i>Sincere</i>	$+(S)\lambda$	λ who believes every announcement he makes
<i>Competent</i>	$+(Co)\lambda$	Everything believed by λ is true
<i>Omniscient</i>	$+(O)\lambda$	λ believes every true formula

2.3 Different Types of Agents in MCL

In MCL, a number of different tags, called *prejudices*, are associated to different types of agents. A tag is defined as $\alpha ::= +(X)\lambda \mid -(X)\lambda$, where X is the name of a prejudice associated to an agent λ . A prejudice is the assumption of a direct dependency between two of three types of logical tokens assertions, beliefs, and facts in the system for one particular agent. For instance, when we say that an agent is sincere, we mean that when he makes an assertion, then he has the belief of that assertion. In detail we introduce the list of the possible combinations in Section 4. We then have the five agent types given in Table 1. The adjective *weak* and *strong* are only applied to the collaborative prejudice (in MCL) and differentiate between an agent who asserts what he believes on *at least one* (i.e., $T_{\diamond}[\lambda : \varphi]$) or *all* (i.e., $T_{\square}[\lambda : \varphi]$) the channels, respectively.

Fig. 1 shows the introduction ($[I.]$) and elimination ($[E.]$) rules for prejudice given in [2]. Note that only negative tags can be derived, i.e., positive tags only appear in the premises of a rule. This is because the intended meaning of a positive prejudice can easily be defined in second order logic quantifying on the formula asserted or believed by an agent. Hence, a positive tag cannot be introduced after one assertion or belief of an agent, e.g., an agent has to know any topic discussed to be tagged as competent $+(Co)$. On the other hand, the existential quantifier characterizing the negative tags allows for the introduction of negative tags, e.g., if an agent does not know a single topic, he is tagged as incompetent $-(Co)$. This is not explicitly formalized in [2] and we introduce a fully-fledged formalization of tags in Section 3. Furthermore, in [2] the weak and strong adjective are used only for a collaborative agent and, e.g., the rules for the sincere agent (rules R.36 and R.41 in Fig. 1) only consider assertions in at least one channel. In our formalization, we consider weak and strong prejudice for each assertion used in a rule.

3 Agents in MAS

In MCL, agents are defined by using three main components of the framework: the sets of announcements, beliefs and facts, i.e., \mathbb{A}_{λ} , \mathbb{B}_{λ} and \mathbb{F} . A natural step is to define how many different types of agents can be defined out of these three sets. To do that, we first extend the results of [2] by considering the relations between these three sets and then use these relations to define agents in MCL.

Intuitively, we can define the following three relations:

$$\begin{array}{ll}
\text{R.32 } \frac{\varphi + (O)\lambda}{\text{B}[\lambda : \varphi]} \quad [\text{E.} + (O)] & \text{R.33 } \frac{\text{B}[\lambda : \varphi] + (Co)\lambda}{\varphi} \quad [\text{E.} + (Co)] \\
\text{R.34 } \frac{\text{B}[\lambda : \varphi] + (Wci)\lambda}{\text{T}_\diamond[\lambda : \varphi]} \quad [\text{E.} + (Wci)] & \text{R.35 } \frac{\text{B}[\lambda : \varphi] + (Sci)\lambda}{\text{T}_\square[\lambda : \varphi]} \quad [\text{E.} + (Sci)] \\
\text{R.36 } \frac{\text{T}_\diamond[\lambda : \varphi] + (S)\lambda}{\text{B}[\lambda : \varphi]} \quad [\text{E.} + (S)] & \text{R.37 } \frac{\sim \text{B}[\lambda : \varphi] \quad \varphi}{-(O)\lambda} \quad [\text{I.} - (O)] \\
\text{R.38 } \frac{\text{B}[\lambda : \varphi] \quad \neg\varphi}{-(Co)\lambda} \quad [\text{I.} - (Co)] & \text{R.39 } \frac{\text{B}[\lambda : \varphi] \quad \sim \text{T}_\diamond[\lambda : \varphi]}{-(Wci)\lambda} \quad [\text{I.} - (Wci)] \\
\text{R.40 } \frac{\text{B}[\lambda : \varphi] \quad \sim \text{T}_\square[\lambda : \varphi]}{-(Sci)\lambda} \quad [\text{I.} - (Sci)] & \text{R.41 } \frac{\text{T}_\diamond[\lambda : \varphi] \quad \sim \text{B}[\lambda : \varphi]}{-(S)\lambda} \quad [\text{I.} - (S)]
\end{array}$$

Fig. 1. The rules for prejudice in MCL

- *Collaboration* $(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$ is the relation between beliefs and announcements of an agent λ . This relation defines the level of collaboration of λ as the *quantity* of data an agent announces with respect to the data he believes. For example, if an agent asserts everything he believes, he is collaborative (recall that belief can be false, in which case the agent might not be competent).
- *Competence* $(\mathbb{B}_\lambda, \mathbb{F})$ is the relation between beliefs of an agent λ and true facts. This relation defines the level of competence of λ and is related to the *quality* of data an agent produces. For example, if everything an agent believes is also true, he is competent (note that this is not the definition of knowledge since an agent could believe in false formulae).
- *Honesty* $(\mathbb{A}_\lambda, \mathbb{F})$ is the relation between announcements made by an agent λ and true facts. This relation defines the level of honesty of λ . For example, if everything an agent shares on a channel is also true, then he is honest.⁴

Given that these three relations are over sets, they express *mereological* relations. We use the *Region Connection Calculus (RCC)* to reason on the different “levels” of collaboration/competence/honesty and to identify which are the different possible relations between the three sets \mathbb{A}_λ , \mathbb{B}_λ and \mathbb{F} . This ultimately defines how many different types of agents we can theoretically consider.

RCC, as defined in [8,4], is an axiomatization of certain spatial concepts and relations in first-order logic. In its broader definition, the RCC theory is composed by eight axioms, and is known as RCC-8, but here we restrict to *RCC-5* by not considering tangential connections between spatial regions. We discuss the choice of RCC-5 in more detail in Section 5.

We define *parthood* as the primitive binary inclusion relation \subseteq , which is reflexive, antisymmetric and transitive. In Table 2, we define the relations of RCC-3, RCC-5 and RCC-8, where X , Y and Z are sets (spatial regions) of formulae and *Connects with* expresses the parthood relation. By applying these relations to the pairs $(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$, $(\mathbb{B}_\lambda, \mathbb{F})$ and $(\mathbb{A}_\lambda, \mathbb{F})$ we can distinguish between different levels of collaboration, competence, and honesty. Every tuple representing the

⁴ Note here that honesty is not necessarily related to correctness. In fact, we define an agent as honest if he asserts the truth even if he does not believe in what he asserts.

Table 2. RCC-3, RCC-5, and RCC-8 relations between spatial regions X, Y and Z

RCC-3 RCC-5 RCC-8	Name	Notation	Definition
	Connects with	$C(X, Y)$	$X \subseteq Y$
	Disconnected from	$\neg C(X, Y)$	$X \not\subseteq Y$
	Part of	$P(X, Y)$	$\forall Z C(Z, X) \rightarrow C(Z, Y)$
	Overlaps	$O(X, Y)$	$\exists Z P(Z, X) \wedge P(Z, Y)$
●	Overlaps Not Equal	$ONE(X, Y)$	$O(X, Y) \wedge \neg EQ(X, Y)$
● ● ●	Equal to	$EQ(X, Y)$	$P(X, Y) \wedge P(Y, X)$
● ● ●	DiscRete from	$DR(X, Y)$	$\neg O(X, Y)$
● ● ●	Partial-Overlap	$PO(X, Y)$	$O(X, Y) \wedge \neg P(X, Y) \wedge \neg P(Y, X)$
● ● ●	Proper-part-of	$PP(X, Y)$	$P(X, Y) \wedge \neg P(Y, X)$
● ● ●	Proper-part-of-inverse	$PPi(X, Y)$	$P(Y, X) \wedge \neg P(X, Y)$
● ● ●	Externally Connected	$EC(X, Y)$	$C(X, Y) \wedge \neg O(X, Y)$
● ● ●	Tangential PP	$TPP(X, Y)$	$PP(X, Y) \wedge \exists Z [EC(Z, X), EC(Z, Y)]$
● ● ●	Tangential PPI	$TPPi(X, Y)$	$TPP(Y, X)$
● ● ●	Non-tangential PP	$NTPP(X, Y)$	$PP(X, Y) \wedge \neg \exists Z [EC(Z, X), EC(Z, Y)]$
● ● ●	Non-tangential PPI	$NTPPi(X, Y)$	$NTPP(Y, X)$

combination of the three relations defines a different type of agent.

$$Agent = \langle RCC5_1(\mathbb{A}_\lambda, \mathbb{B}_\lambda), RCC5_2(\mathbb{B}_\lambda, \mathbb{F}), RCC5_3(\mathbb{F}, \mathbb{A}_\lambda) \rangle$$

where $RCC5_1$, $RCC5_2$ and $RCC5_3$ are relations in RCC-5. As we discuss in Section 5, some combinations of $RCC5_1$, $RCC5_2$ and $RCC5_3$ are topologically incorrect.

4 Categorization of Agents

We now consider the details of every RCC-5 relation between each pair of \mathbb{A}_λ , \mathbb{B}_λ and \mathbb{F} and we define 15 different prejudices. Our list is complete with respect to RCC-5, i.e., no other relations can be considered. We will use roman numerals to identify the new rules we introduce, whereas the decimals for the rules were already defined in [2].

4.1 Collaboration

Sincere $PP(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$. A sincere agent λ is defined by the proper part of his announcements with respect to his beliefs. More formally, for any propositional formula φ ,

$$\text{if } T_*[\lambda : \varphi] \text{ then } B[\lambda : \varphi],$$

where $*$ identifies one of the two modalities in MCL, i.e., $*$ \in $\{\Box, \Diamond\}$.

This type of agent announces *only* what he believes (\Rightarrow) but does not announce everything he believes (\neq). As already defined in [2], we can negate the

formula of a sincere agent and provide deduction rules to define when an agent is *not* sincere as follows. For a non-sincere agent λ^5 , there exists a propositional formula φ such that

$$T_*[\lambda : \varphi] \text{ and } \sim B[\lambda : \varphi].$$

We can then define rules that formalize that if an agent asserts, even only once, something that he does not believe in, then he is non-sincere:

$$\text{R.41 } \frac{T_\diamond[\lambda : \varphi] \quad \sim B[\lambda : \varphi]}{-(W_S)\lambda} \quad [I.-(W_S)] \quad \text{R.I } \frac{T_\square[\lambda : \varphi] \quad \sim B[\lambda : \varphi]}{-(S_S)\lambda} \quad [I.-(S_S)]$$

As we discussed in Section 2.3, the notion of weak and strong is only applied to the notion of collaborative agent in MCL. We avoid this asymmetry and we introduce the notion of weak and strong for all the prejudices involving a relation with announcements. This explains why we have now used W_S in R.41 instead of S of MCL (as in Fig. 1). We extend the elimination rules accordingly:

$$\begin{array}{ll} \text{R.II } \frac{\sim B[\lambda : \varphi] \quad + (W_S)\lambda}{\sim T_\diamond[\lambda : \varphi]} \quad [E.+(W_S)] & \text{R.III } \frac{\sim B[\lambda : \varphi] \quad + (S_S)\lambda}{\sim T_\square[\lambda : \varphi]} \quad [E.+(S_S)] \\ \text{R.36 } \frac{T_\diamond[\lambda : \varphi] \quad + (W_S)\lambda}{B[\lambda : \varphi]} \quad [E.+(W_S)] & \text{R.IV } \frac{T_\square[\lambda : \varphi] \quad + (S_S)\lambda}{B[\lambda : \varphi]} \quad [E.+(S_S)] \end{array}$$

Collaborative $PPi(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$. Symmetrically to a sincere agent, a collaborative agent λ is defined by the proper part of his beliefs with respect to his announcements: for any propositional formula φ ,

$$\text{if } B[\lambda : \varphi] \text{ then } T_*[\lambda : \varphi].$$

This type of agent announces everything he believes (\Rightarrow) but what he says is *not only* what he believes (\neq). Hence, some of the announcements are intentionally against his beliefs (these announcements might be accidentally true facts but we will discuss this case later in this section). If we negate the definition of collaborative, we obtain that if an λ 's belief has not been announced (i.e., there exists φ such that $B[\lambda : \varphi]$ and $\sim T_*[\varphi : \lambda]$), then λ is *not* collaborative. As for the sincere agent, we define strong and weak prejudice with \square and \diamond , respectively:

$$\begin{array}{ll} \text{R.39 } \frac{B[\lambda : \varphi] \quad \sim T_\diamond[\lambda : \varphi]}{-(W_{CI})\lambda} \quad [I.-(W_{CI})] & \text{R.40 } \frac{B[\lambda : \varphi] \quad \sim T_\square[\lambda : \varphi]}{-(S_{CI})\lambda} \quad [I.-(S_{CI})] \\ \text{R.34 } \frac{B[\lambda : \varphi] \quad + (W_{CI})\lambda}{T_\diamond[\lambda : \varphi]} \quad [E.+(W_{CI})] & \text{R.35 } \frac{B[\lambda : \varphi] \quad + (S_{CI})\lambda}{T_\square[\lambda : \varphi]} \quad [E.+(S_{CI})] \\ \text{R.V } \frac{\sim T_\diamond[\lambda : \varphi] \quad + (W_{CI})\lambda}{\sim B[\lambda : \varphi]} \quad [E.+(W_{CI})] & \text{R.VI } \frac{\sim T_\square[\lambda : \varphi] \quad + (S_{CI})\lambda}{\sim B[\lambda : \varphi]} \quad [E.+(S_{CI})] \end{array}$$

⁵ Slightly abusing notation, we are using λ both for a sincere and non-sincere agent.

Fair $EQ(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$. A fair agent λ is defined by the equality between the sets of his announcements and beliefs: for any propositional formula φ ,

$$T_*[\lambda : \varphi] \text{ if and only if } B[\lambda : \varphi].$$

Hence, a fair agent is an agent who believes in *everything* he announces (\Rightarrow) and who announces *only* what he believes (\Leftarrow). As before, in order to give the rules for MCL, we first negate the definition of the fair agent. For a non-fair agent λ , there exists a propositional formula φ such that

$$(\sim T_*[\lambda : \varphi] \text{ and } B[\lambda : \varphi]) \text{ or } (\sim B[\lambda : \varphi] \text{ and } T_*[\lambda : \varphi]).$$

The left and right disjuncts are exactly the definitions of PPI and PP, respectively. Hence, the introduction and elimination rules have been already considered in the previous two cases.

Saboteur $PO(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$. A saboteur agent λ is defined by the partial overlap of his announcements with respect to his beliefs: for any propositional formula φ ,

$$B[\lambda : \varphi] \text{ or } T_*[\lambda : \varphi].$$

This type of agent may announce something that he believes but also that he does not believe, or does not announce something he believes.

$$\text{R.VII } \frac{\sim B[\lambda : \varphi] \quad \sim T_\diamond[\lambda : \varphi]}{-(W_I)\lambda} \quad [\text{I.}-(W_I)] \quad \text{R.VIII } \frac{\sim B[\lambda : \varphi] \quad \sim T_\square[\lambda : \varphi]}{-(S_I)\lambda} \quad [\text{I.}-(S_I)]$$

$$\text{R.IX } \frac{\sim T_\diamond[\lambda : \varphi] \quad +(S_I)\lambda}{B[\lambda : \varphi]} \quad [\text{E.}+(S_I)] \quad \text{R.X } \frac{\sim T_\square[\lambda : \varphi] \quad +(S_I)\lambda}{B[\lambda : \varphi]} \quad [\text{E.}+(S_I)]$$

$$\text{R.XI } \frac{\sim B[\lambda : \varphi] \quad +(W_I)\lambda}{T_\diamond[\lambda : \varphi]} \quad [\text{E.}+(W_I)] \quad \text{R.XII } \frac{\sim B[\lambda : \varphi] \quad +(W_I)\lambda}{T_\square[\lambda : \varphi]} \quad [\text{E.}+(W_I)]$$

Braggart $DR(\mathbb{A}_\lambda, \mathbb{B}_\lambda)$. A braggart agent λ is defined by the discrete-from relation between his announcements and beliefs: for any propositional formula φ ,

$$\sim T_*[\lambda : \varphi] \text{ or } \sim B[\lambda : \varphi].$$

This agent *only* announces what he does not believe and he does not announce what he believes. Reasoning on the negated definition (i.e., on a non-braggart agent λ for which there exists a propositional formula φ such that $T_*[\lambda : \varphi]$ and $B[\lambda : \varphi]$), we can define that if (at least once) the agent states something he believes in, then he is non-braggart.

$$\text{R.XIII } \frac{T_\diamond[\lambda : \varphi] \quad B[\lambda : \varphi]}{-(W_B)\lambda} \quad [\text{I.}-(W_B)] \quad \text{R.XIV } \frac{T_\square[\lambda : \varphi] \quad B[\lambda : \varphi]}{-(S_B)\lambda} \quad [\text{I.}-(S_B)]$$

$$\text{R.XV } \frac{T_\diamond[\lambda : \varphi] \quad +(S_B)\lambda}{\sim B[\lambda : \varphi]} \quad [\text{E.}+(S_B)] \quad \text{R.XVI } \frac{T_\square[\lambda : \varphi] \quad +(S_B)\lambda}{\sim B[\lambda : \varphi]} \quad [\text{E.}+(S_B)]$$

$$\text{R.XVII } \frac{B[\lambda : \varphi] \quad +(W_B)\lambda}{\sim T_\diamond[\lambda : \varphi]} \quad [\text{E.}+(W_B)] \quad \text{R.XVIII } \frac{B[\lambda : \varphi] \quad +(W_B)\lambda}{\sim T_\square[\lambda : \varphi]} \quad [\text{E.}+(W_B)]$$

4.2 Competence

Competent $PP(\mathbb{B}_\lambda, \mathbb{F})$. An agent's beliefs are a subset of the true formulae. Hence, all the agent's beliefs are facts but there may be true formulae "out" of his beliefs. An agent λ is competent if, for every propositional formula φ , if $B[\lambda : \varphi]$ then $\varphi \in \mathbb{F}$.

$$\text{R.38} \frac{B[\lambda : \varphi] \quad \neg\varphi}{-(Co)\lambda} \quad [I. -(Co)]$$

$$\text{R.33} \frac{B[\lambda : \varphi] \quad +(Co)\lambda}{\varphi} \quad [E. +(Co)] \quad \text{R.XIX} \frac{\sim B[\lambda : \varphi] \quad +(Co)\lambda}{\neg\varphi} \quad [E. +(Co)]$$

Omniscient $PP_i(\mathbb{B}_\lambda, \mathbb{F})$. An agent λ is omniscient if the set of formulae he believes is a superset of the actually true formulae: for any propositional formula φ , if $\varphi \in \mathbb{F}$ then $B[\lambda : \varphi]$.

$$\text{R.37} \frac{\sim B[\lambda : \varphi] \quad \varphi}{-(O)\lambda} \quad [I. -(O)]$$

$$\text{R.32} \frac{\varphi \quad +(O)\lambda}{B[\lambda : \varphi]} \quad [E. +(O)] \quad \text{R.XX} \frac{\sim B[\lambda : \varphi] \quad +(O)\lambda}{\neg\varphi} \quad [E. +(O)]$$

Wise $EQ(\mathbb{B}_\lambda, \mathbb{F})$. A wise agent λ is defined by the equality between the sets of his beliefs and facts, i.e., he *only* believes in true formulae and knows *all* the true facts: for any propositional formula φ , $\varphi \in \mathbb{F}$ if and only if $B[\lambda : \varphi]$. The rules generated are exactly the rules of PP_i and PP .

Incompetent $PO(\mathbb{B}_\lambda, \mathbb{F})$. An incompetent agent λ is defined by the partial overlap of his beliefs with the true facts: for any propositional formula φ , $\varphi \in \mathbb{F}$ or $B[\lambda : \varphi]$. This type of agent believes in true and false formulae, and there exist facts that he does not believe in, but he won't believe a false formula φ .

$$\text{R.XXI} \frac{\neg\varphi \quad \sim B[\lambda : \varphi]}{-(In)\lambda} \quad [I. -(In)]$$

$$\text{R.XXII} \frac{\neg\varphi \quad +(In)\lambda}{B[\lambda : \varphi]} \quad [E. +(In)] \quad \text{R.XXIII} \frac{\sim B[\lambda : \varphi] \quad +(In)\lambda}{\varphi} \quad [E. +(In)]$$

Ignorant $DR(\mathbb{B}_\lambda, \mathbb{F})$. An ignorant agent λ is defined by the discrete-from relation between true formulae and beliefs: for any propositional formula φ , $\neg\varphi \in \mathbb{F}$ or $\sim B[\lambda : \varphi]$. Therefore, this agent *only* believes in false formulae.

$$\text{R.XXIV} \frac{\varphi \quad B[\lambda : \varphi]}{-(Ig)\lambda} \quad [I. -(Ig)]$$

$$\text{R.XXV} \frac{\varphi \quad +(Ig)\lambda}{\sim B[\lambda : \varphi]} \quad [E. +(Ig)] \quad \text{R.XXVI} \frac{B[\lambda : \varphi] \quad +(Ig)\lambda}{\neg\varphi} \quad [E. +(Ig)]$$

4.3 Honesty

Honest $PP(\mathbb{A}_\lambda, \mathbb{F})$. An agent is honest if every formula he asserts is a fact, and the agent's assertion are a subset of the true formulae: for any propositional formula φ , if φ then $T_*[\lambda : \varphi]$.

$$\begin{array}{ll}
\text{R.XXVII} \frac{\varphi \sim T_\diamond[\lambda : \varphi]}{-(W_H)\lambda} \quad [\text{I. } -(W_H)] & \text{R.XXVIII} \frac{\varphi \sim T_\square[\lambda : \varphi]}{-(S_H)\lambda} \quad [\text{I. } -(S_H)] \\
\text{R.XXIX} \frac{\varphi + (W_H)\lambda}{T_\diamond[\lambda : \varphi]} \quad [\text{E. } +(W_H)] & \text{R.XXX} \frac{\varphi + (S_H)\lambda}{T_\square[\lambda : \varphi]} \quad [\text{E. } +(S_H)] \\
\text{R.XXXI} \frac{\sim T_\diamond[\lambda : \varphi] + (W_H)\lambda}{\neg\varphi} \quad [\text{E. } +(W_H)] & \text{R.XXXII} \frac{\sim T_\square[\lambda : \varphi] + (S_H)\lambda}{\neg\varphi} \quad [\text{E. } +(S_H)]
\end{array}$$

Oracle $PP_i(\mathbb{A}_\lambda, \mathbb{F})$. An agent λ is an oracle if, for any propositional formula φ , if $T_*[\lambda : \varphi]$ then $\varphi \in \mathbb{F}$.

$$\begin{array}{ll}
\text{R.XXXIII} \frac{T_\diamond[\lambda : \varphi] \neg\varphi}{X} \quad [\text{I. } -(W_{Or})] & \text{R.XXXIV} \frac{T_\square[\lambda : \varphi] \neg\varphi}{X} \quad [\text{I. } -(S_{Or})] \\
\text{R.XXXV} \frac{T_\diamond[\lambda : \varphi] + (W_{Or})\lambda}{\varphi} \quad [\text{E. } +(W_{Or})] & \text{R.XXXVI} \frac{T_\square[\lambda : \varphi] + (S_{Or})\lambda}{\varphi} \quad [\text{E. } +(S_{Or})] \\
\text{R.XXXVII} \frac{\neg\varphi + (W_{Or})\lambda}{\sim T_\diamond[\lambda : \varphi]} \quad [\text{E. } +(W_{Or})] & \text{R.XXXVIII} \frac{\neg\varphi + (S_{Or})\lambda}{\sim T_\square[\lambda : \varphi]} \quad [\text{E. } +(S_{Or})]
\end{array}$$

Right $EQ(\mathbb{A}_\lambda, \mathbb{F})$. An agent λ is right if, for any propositional formula φ , $\varphi \in \mathbb{F}$ if and only if $T_*[\lambda : \varphi]$. We omit the rules since they are the same as for PP and PP_i .

Incorrect $PO(\mathbb{A}_\lambda, \mathbb{F})$. An agent λ is incorrect if, for any propositional formula φ , $\varphi \in \mathbb{F}$ or $T_*[\lambda : \varphi]$. The announcements of this type of agent might be true or false, and he only announces part of the facts (i.e., a subset of the facts will never be announced by him).

$$\begin{array}{ll}
\text{R.XXXIX} \frac{\neg\varphi \sim T_\diamond[\lambda : \varphi]}{-(W_{Ir})\lambda} \quad [\text{I. } -(W_{Ir})] & \text{R.XL} \frac{\neg\varphi \sim T_\square[\lambda : \varphi]}{-(S_{Ir})\lambda} \quad [\text{I. } -(S_{Ir})] \\
\text{R.XLI} \frac{\neg\varphi + (W_{Ir})\lambda}{T_\diamond[\lambda : \varphi]} \quad [\text{E. } +(W_{Ir})] & \text{R.XLII} \frac{\neg\varphi + (S_{Ir})\lambda}{T_\square[\lambda : \varphi]} \quad [\text{E. } +(S_{Ir})] \\
\text{R.XLIII} \frac{\sim T_\diamond[\lambda : \varphi] + (W_{Ir})\lambda}{\varphi} \quad [\text{E. } +(W_{Ir})] & \text{R.XLIV} \frac{\sim T_\square[\lambda : \varphi] + (S_{Ir})\lambda}{\varphi} \quad [\text{E. } +(S_{Ir})]
\end{array}$$

False $DR(\mathbb{A}_\lambda, \mathbb{F})$. A false agent λ is defined by the discrete-form relation between true formulae and his assertions, i.e., for any propositional formula φ , $\neg\varphi \in \mathbb{F}$ or $\sim T_*[\lambda : \varphi]$. In other words, everything he announces is false.

$$\begin{array}{ll}
\text{R.XLV} \frac{\varphi \ T_{\diamond}[\lambda : \varphi]}{-(W_F)\lambda} \quad [\text{I.}-(W_F)] & \text{R.XLVI} \frac{\varphi \ T_{\square}[\lambda : \varphi]}{-(S_F)\lambda} \quad [\text{I.}-(S_F)] \\
\text{R.XLVII} \frac{\varphi \ + (W_F)\lambda}{\sim T_{\diamond}[\lambda : \varphi]} \quad [\text{E.}+(W_F)] & \text{R.XLVIII} \frac{\varphi \ + (S_F)\lambda}{\sim T_{\square}[\lambda : \varphi]} \quad [\text{E.}+(S_F)] \\
\text{R.XLIX} \frac{T_{\diamond}[\lambda : \varphi] \ + (W_F)\lambda}{\neg\varphi} \quad [\text{E.}+(W_F)] & \text{R.L} \frac{T_{\square}[\lambda : \varphi] \ + (S_F)\lambda}{\neg\varphi} \quad [\text{E.}+(S_F)]
\end{array}$$

5 On the Topology of MAS

In this section, we justify the use of RCC-5 instead of RCC-3 or RCC-8, and discuss the relation between the topology we consider and the agent types.

5.1 RCC-3, RCC-5, and RCC-8

As already mentioned in Section 3, there exist three different types of RCC, based on the number of topological relations considered: RCC-3, RCC-5, and RCC-8. RCC-3 considers the three different topological relations listed in Table 2: *ONE*, *EQ*, and *DR*. The topological relations *EQ* and *DR* are the same as in RCC-5 (see Table 2), whereas *ONE* defines the overlap relation between two regions with the additional constraint that the regions cannot be fully overlapping (i.e., they cannot be two exact copies of the same region).

The relation *ONE* in RCC-3 is detailed in RCC-5 with the relations *PP*, *PPi*, and *PO*. Hence, considering RCC-5 instead of RCC-3 results in a more accurate and expressive categorization of agents. However, the same reasoning cannot be applied to RCC-8. In fact, even if RCC-8 is more detailed than RCC-5 as it considers more topological relations, the additional topological relations considered by RCC-8 cannot be applied for the categorization of agents in MCL. As showed in Table 2, RCC-8 considers tangential connections, where, informally, two tangential regions are near enough so that no other region can fit between the two (without overlapping them), but are not overlapping at any point. This is formalized by the *EC* relation. In addition, in RCC-8, each of the two relations *PP* and *PPi* is detailed into tangential and non-tangential.

In our work, the elements of the three sets \mathbb{A} , \mathbb{B} and \mathbb{F} are not ordered. In other words, we are not considering the distance between those elements (or between regions containing those elements). Hence, given any pair of (sub-)sets between \mathbb{A} , \mathbb{B} and \mathbb{F} , regardless of the sets being near or far apart between each other, we consider them as disjoint (i.e., *DR*).

5.2 An Upper-bound on the Number of Different Types of Agents

Applying RCC over a finite number of sets, we obtain a definite number of resulting combinations. Hence, applying RCC over $\mathbb{A}_\lambda, \mathbb{B}_\lambda, \mathbb{F}$, we obtain a definite number of different types of agents. In this section, we show the general upper

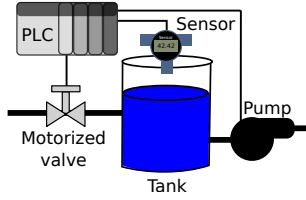


Fig. 2. Representation of the test case

	Theoretical	Correct
<i>RCC-3</i>	$3^3 = 27$	15
<i>RCC-5</i>	$5^3 = 125$	54
<i>RCC-8</i>	$8^3 = 512$	193

Fig. 3. Number of agents with respect to different RCC

bound on the number of different agents with respect to the type of RCC (RCC-5, RCC-3 or RCC-8) considered.

The general formula to calculate the number of different types of agents is $r^{\binom{n}{k}}$, where r is the number of relations with arity k , between n different sets, where r^e is the number of permutation of r relations over e elements with repetitions, with e being the number of k -ary combinations of n sets, $\binom{n}{k}$. In our case, $\binom{3}{k} = 3$ since we consider 3 sets ($\mathbb{A}, \mathbb{B}, \mathbb{F}$), and all the relations considered in the RCC are binary. Hence, using RCC-5 (with five different spatial relations) over three sets, we can theoretically define up to 125 different type of agents. However, only 54 of the 125 (as showed in [4] and derived by the composition table of RCC-5) combinations are topologically correct with respect to the definition of the relations of RCC-5. Generalizing to all the RCCs, in Tab. 3 we calculate the number of different agents with respect to all the variations of RCC (i.e., with 3, 5 or 8 spatial relations). Due to space limits, we omit the composition table for RCC-3, RCC-5 and RCC-8. Hence, even if considering a different number of sets than the three \mathbb{A} , \mathbb{B} and \mathbb{F} exponentially affects the number of theoretical agents, the application of RCC downscales that number of a factor that ranges from 1.8 to 2.5. In addition, using RCC-5 we consider 3.6 times more (different) types of agents than RCC-3, but using RCC-8 would allow us to consider 3.5 times more different agents.

6 Use Case

In this section, we show that both the framework and the categorization of agents that we have given can be applied to reason about the security of CPS.

6.1 Cyber-Physical Systems

We use the term CPS to refer to systems that consist of networked embedded systems, which are used to sense, actuate, and control physical processes. Examples of CPS include industrial water treatment facilities and power plants. CPS have seen a rapid increase in automation and connectivity, which threatens to increase their vulnerability to malicious attacks. Let us now use our approach to address the problem of defining security-related attack states for CPS.

Description of the Case Study. Similarly to [9,5] we consider a CPS (depicted in Figure 2) composed by five agents:

Table 3. Example of attack states for the water level sensor

State of the sensor	(A, B)	(B, F)	(A, F)
optimal	EQ	EQ	EQ
sensor compromised	EQ	DR	DR
communication compromised	DR	EQ	DR
fully compromised	DR	DR	DR

- A *tank* containing water.
- A *controller* (e.g., a PLC) that controls the water level so that the tank does not (underflow or) overflow.
- A *water level indicator* (e.g., a Sensor) that communicates the readings of the level of the water inside the tank to the PLC.
- A *motorized valve* and a *pump* that (controlled by the PLC) regulate the inflow and outflow of water respectively.

Mapping \mathbb{A} , \mathbb{B} , and \mathbb{F} to CPS. It is possible that the three sets \mathbb{A} , \mathbb{B} and \mathbb{F} contain at the same time different formulae that contain each element of the topological space φ . Hence, every assertion and belief must be objective (since it can be part of \mathbb{F}). This implies that formulae like $\varphi := \text{highLevel}(\text{tank}, \text{water})$ cannot be considered in our reasoning since “high” is considered to be subjective. In contrast, we can use objective formulae such as $\varphi := \text{level}(\text{tank}, 20L)$.

When considering a CPS (and security systems in general, e.g., security protocols) as a MAS, the message exchange between different agents can be formalized by means of *assertions*. In addition, redundant channels are often employed to reduce security treats (or assertions are required over multiple channels as, e.g., in two-factor authentication) and then it is fair to assume that assertions can be done over single or multiple channels. Finally, the inspection of the memory of any software/hardware of the CPS (supposing a white-box analysis) reveals the actual *beliefs*, while the *facts* in a CPS are defined by the physical laws of the physics. We can summarize our mapping as follows:

- \mathbb{A}_λ defines the values communicated by the agent λ .
- \mathbb{B}_λ defines the computational results of the agent λ .
- \mathbb{F} defines the environmental values, i.e., the real values of the system.

6.2 Single-Channel Attack states

We are now in a position to show that we can directly apply our topological categorization to any agent in our CPS. For simplicity, we first use only the RCC-5 relations EQ and DR, and then extend our results to all RCC-5 relations.

Optimal System Status. Suppose that the tank contains 20L of water, e.g., $\text{level}(\text{tank}, 20L) \in \mathbb{F}$, where *level* is a predicate, and *tank* and *20L* are propositional constants. For the sake of simplicity, we also suppose that the system is in idle (both the motorized valve and the pump are off). When the system is *not* compromised, the sensor correctly computes the level of the water in the tank (e.g., $\text{level}(\text{tank}, 20L) \in \mathbb{B}_{\text{sensor}}$) and correctly communicates

to the PLC the computed value of water in the tank (e.g., $level(tank, 20L) \in \mathbb{A}_{sensor}$). We can then define the optimal status of the sensor as the triple $\langle EQ(\mathbb{A}_{sensor}, \mathbb{B}_{sensor}), EQ(\mathbb{B}_{sensor}, \mathbb{F}), EQ(\mathbb{A}_{sensor}, \mathbb{F}) \rangle$.

System Under Attack. Suppose that the sensor is communicating wrong values to the PLC (i.e., $DR(\mathbb{A}_{sensor}, \mathbb{F})$). As showed in Table 3, we have three mutually exclusive cases:

1. The sensor is working properly $EQ(\mathbb{B}_{sensor}, \mathbb{F})$, therefore (topologically) the communication between the sensor and the PLC has been compromised, i.e., $DR(\mathbb{A}_{sensor}, \mathbb{B}_{sensor})$.
2. The communication between the sensor and the PLC has not been compromised $EQ(\mathbb{A}_{sensor}, \mathbb{B}_{sensor})$, therefore the sensor is *not* sending what it computes $DR(\mathbb{B}_{sensor}, \mathbb{F})$.
3. Both the communication and the sensor have been compromised.

As a consequence of the discussion in Section 5.2, between the optimal and the fully compromised status of the sensor there must be 52 other different statuses. Due to lack of space, we cannot go into the details of each status, but we can generalize the attack states into three main categories, as follows:

- $RCC5(\mathbb{A}, \mathbb{B})$ expresses the relation between the values communicated and the ones computed by an agent.
- $RCC5(\mathbb{B}, \mathbb{F})$ expresses the relation between the values computed and the true environmental values.
- $RCC5(\mathbb{A}, \mathbb{F})$ expresses the relation between the values communicated and the true environmental values.

Defense mechanisms that check sudden changes in physical readings (for an example of how this is defined in MAS with logical systems, see [12]) are often adopted in CPS. To bypass the security mechanisms, during an attack, the optimal status will likely pass through most of the 52 intermediate statuses.

6.3 Multiple-Channel Attack States

A countermeasure often applied in CPS (but not limited to CPS) is the implementation of redundant channels. As proposed in [10], in our case study one could implement a dedicated system that interprets the readings of the sensor and directly closes the motorized valve if an upper threshold is reached. We can leverage the modal operators to define such communications and to define even more sophisticated attack states. For example, given a state $\mathbb{A}_{sensor}, \mathbb{B}_{sensor}, \mathbb{F}$ in MCL, we can check if one or all the channels that the sensor uses to communicate with the PLC have been compromised, as defined in (1) and (2) respectively:

$$\{\mathbb{A}_{sensor}, \mathbb{B}_{sensor}, \mathbb{F}\} \vdash \neg (S_{Fair})_{sensor} \quad (1)$$

$$\{\mathbb{A}_{sensor}, \mathbb{B}_{sensor}, \mathbb{F}\} \vdash \neg (W_{Fair})_{sensor} \quad (2)$$

Based on the approach we have proposed in this paper we can formalize the optimal/attack states of a CPS, reason on the properties of the CPS by means

of prejudices in MCL, and obtain therefore a control upon the concept of redundancy as expressed above. Our approach is not specific to CPS but can potentially be applied to any MAS (as long as the elements of the topological space are objective).

7 Conclusion

We proposed a topological categorization of agents for MCL using RCC5. We defined an upper bound on the number of different agents in a MAS and we applied our results to the security of CPS. We showed that our results can be used to address the problem of defining attack states for CPS. We are currently working on an implementation of our framework. We have also been extending MCL to capture the intents of agents, which will ultimately allow us to consider human agents in the formalization of MAS.

References

1. P. Balbiani and P. Seban. Reasoning about permitted announcements. *Journal of Philosophical Logic*, 40(4):445–472, 2011.
2. M. Cristani, F. Olivieri, and K. Santacà. A logical model of communication channels. In *Intelligent and Evolutionary Systems*, 2016.
3. D. Grossi, L. Royakkers, and F. Dignum. Organizational structure and responsibility : An analysis in a dynamic logic of organized collective agency. *Artificial Intelligence and Law*, 15(3):223–249, 2007.
4. R. Grütter, T. Scharrenbach, and B. Bauer-Messmer. Improving an rcc-derived geospatial approximation by OWL axioms. In *ISWC*, 2008.
5. E. Kang, S. Adepu, D. Jackson, and A. P. Mathur. Model-based security analysis of a water treatment system. In *SEsCPS*, 2016.
6. S. Khaitan and J. McCalley. Design techniques and applications of cyberphysical systems: A survey. *IEEE Systems Journal*, 9(2):350–365, 2015.
7. J. Lin, S. Sedigh, and A. Miller. Modeling Cyber-Physical Systems with Semantic Agents. In *COMPSACW*, 2010.
8. T. Y. Lin, Q. Liu, and Y. Y. Yao. Logics systems for approximate reasoning: Approximation via rough sets and topological spaces. In *ISMIS*, 1994.
9. M. Rocchetto and N. O. Tippenhauer. CPDY: Extending the Dolev-Yao Attacker with Physical-Layer Interactions. In *ICFEM*, 2016.
10. G. Sabaliauskaite and A. P. Mathur. Intelligent checkers to improve attack detection in cyber physical systems. In *CyberC*, 2013.
11. T. Sanislav and L. Miclea. Cyber-physical systems - concept, challenges and research areas. *Control Engineering and Applied Informatics*, 14(2):28–33, 2012.
12. D. Urbina, J. Giraldo, A. A. Cardenas, N. O. Tippenhauer, J. Valente, M. Faisal, J. Ruths, R. Candell, and H. Sandberg. Limiting the impact of stealthy attacks on industrial control systems. In *CCS*, 2016.
13. J. Van Benthem, J. Van Eijck, and B. Kooi. Logics of communication and change. *Information and computation*, 204(11):1620–1662, 2006.