



ALMA MATER STUDIORUM  
UNIVERSITÀ DI BOLOGNA

ARCHIVIO ISTITUZIONALE  
DELLA RICERCA

## Alma Mater Studiorum Università di Bologna Archivio istituzionale della ricerca

Towards Explainable Visionary Agents: License to Dare and Imagine

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

*Published Version:*

Giovanni Ciatto, A.N. (2021). Towards Explainable Visionary Agents: License to Dare and Imagine. Cham : Springer Nature [10.1007/978-3-030-82017-6\_9].

*Availability:*

This version is available at: <https://hdl.handle.net/11585/842452> since: 2021-12-20

*Published:*

DOI: [http://doi.org/10.1007/978-3-030-82017-6\\_9](http://doi.org/10.1007/978-3-030-82017-6_9)

*Terms of use:*

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).  
When citing, please refer to the published version.

(Article begins on next page)

This is the final peer-reviewed accepted manuscript of:

**Ciatto, G., Najjar, A., Calbimonte, JP., Calvaresi, D. (2021). Towards Explainable Visionary Agents: License to Dare and Imagine. In: Calvaresi, D., Najjar, A., Winikoff, M., Främling, K. (eds) Explainable and Transparent AI and Multi-Agent Systems. EXTRAAMAS 2021. Lecture Notes in Computer Science(), vol 12688. Springer, Cham.**

The final published version is available online at: [https://doi.org/10.1007/978-3-030-82017-6\\_9](https://doi.org/10.1007/978-3-030-82017-6_9)

Rights / License:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

*This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>)*

***When citing, please refer to the published version.***

# Towards Explainable Visionary Agents: License to Dare and Imagine

Giovanni Ciatto <sup>✉</sup><sup>2</sup>[0000-0002-1841-8996], Amro Najjar<sup>3</sup>[0000-0001-7784-6176],  
Jean-Paul Calbimonte<sup>1</sup>[0000-0001-7784-6176], and Davide  
Calvaresi<sup>1</sup>[0000-0001-9816-7439]

<sup>1</sup> University of Applied Sciences and Arts Western Switzerland HES-SO, Switzerland

{davide.calvaresi, jean-paul.calbimonte}@hevs.ch

<sup>2</sup> University of Bologna, Bologna, Italy

giovanni.ciatto@unibo.it

<sup>3</sup> University of Luxembourg, Luxembourg

amro.najjar@uni.lu

**Abstract.** Since their appearance, computer programs have embodied discipline and structured approaches and methodologies. Yet, to this day, equipping machines with imaginative and creative capabilities remains one of the most challenging and fascinating goals we pursue. Intelligent software agents can behave *intelligently* in well-defined scenarios, relying on Machine Learning (ML), symbolic reasoning, and the ability of their developers for tailoring *smart behaviors* to specific application domains. However, to forecast the evolution of all possible scenarios is unfeasible. Thus, intelligent agents should autonomously/creatively adapt to the world's mutability. This paper investigates the meaning of imagination in the context of cognitive agents. In particular, it addresses techniques and approaches to let agents *autonomously imagine/simulate* their course of action and generate explanations supporting it, and formalizes thematic challenges. Accordingly, we investigate research areas including: *(i)* reasoning and automatic theorem proving to synthesize novel knowledge via inference; *(ii)* automatic planning and simulation, used to speculate over alternative courses of action; *(iii)* machine learning and data mining, exploited to induce new knowledge from experience; and *(iv)* biochemical coordination, which keeps imagination dynamic by continuously reorganizing it.

**Keywords:** Multi-agent systems, imagination, BDI, cognitive agents, XAI.

## 1 Introduction

Imagination is among the most powerful tools humankind has ever had. Fundamentally, imagination is responsible for the *spontaneous* creation of novel ideas which do not originate directly from the human senses. Such a mental process enabled humans to design complex concepts and artifacts, shaping the societies we live in nowadays [1].

Over the years, imagination and creativity have been considered the nemesis of discipline and structured approaches in general [2]. Indeed, at an individual level, imagination is a relatively simple process. It can be conceived as *a never-ending activity occurring within each person’s mind along their entire lifetime*. Such activity copes with reorganizing a person’s beliefs, perceptions, feelings, and habits repeatedly to *generate* novel beliefs, abilities, desires, insights about the future, and needs—which in turn may motivate novel activities.

Let us consider two simple examples commonly dealt with:

**Counterfactual thinking:** it is a retrospective “what if” analysis, elaborating how things could have been different (e.g., regretting a decision, “*I should have behaved differently*”) – also known as staircase wit – from which a *lesson* is supposably learned.

**Speculative thinking:** mentally simulating possible future scenarios according to models of (i) the world, and (ii) other agents/humans behaviours—e.g., imaging the effect a proposed example might have on the reader.

There, imagination is a key enabler for intelligent behavior.

In modern Artificial Intelligence (AI), many research efforts are devoted to the engineering of *smart mechanisms*, which enable software agents to behave *intelligently* in well-defined scenarios. Most of these mechanisms are either based on Machine Learning (ML) or on symbolic reasoning (including planning or automatic theorem proving) [3]. Nevertheless, the capability of software agents (intended as intelligent virtual entities) to behave intelligently strongly depends on their developers’ capability of tailoring *smart behaviors* on the particular scenario the software agents operate into. Arguably, however, it is unfeasible for developers to forecast all possible evolutions a real-world scenario may be subject to. Accordingly, intelligent software agents should also adapt to the world by autonomously figuring out how to deal with its mutability—similarly to what a human would do. Notably, one of the significant areas where adaptability is expected to play a major role is XAI. Indeed, there is an increasing push for intelligent systems capable of explaining their own behavior [4]. However, current research efforts are mostly focused on supporting data scientists in drawing explanations in particular cases. Even when explanations are delegated to autonomous agents, their capability to generate effective explanations still relies on their developers’ foresight. In other words, the problem of letting agents *autonomously* generate explanations is still open.

Human beings heavily leverage on imagination to adapt to the world. In particular, they exploit both counterfactual and speculative thinking to adapt the way they *interact* with their interlocutor. Arguably, similar mechanisms could be conceived for agents willing to attain the capability of *generating* explanations.

Accordingly, in this paper, we discuss (i) what imagination may mean for software agents, (ii) how it could be technically realized within modern agent frameworks, (iii) what is the role of imagination-equipped agents in modern data-driven AI, and (iv) how can imagination support the autonomous generation of explanations. In doing so, we restrict our scope to the case of *cognitive* agents, as their abstractions are rich enough to capture a general – yet precise

– notion of imagination. In particular, we focus on the Belief-Desire-Intention (BDI) agent architectures as they represent the best viable bridge among theory and practice—being backed by effective technologies such as Jason [5].

Within the scope of this paper, we conceive *imagination* as a non-terminating background activity carried on by an agent behind the scenes, possibly while doing anything else. The imagination activity takes care of *continuously* revising an agent’s internal knowledge, possibly (i) obliterating useless information; (ii) synthesizing novel information out of the current and previous experience; (iii) dismissing or generating desires and needs; (iv) critically analyzing the previous courses of actions w.r.t. their goals; (v) simulating possible similar/alternative behaviors to be exploited in similar situations; and (vi) looking for *post-hoc* motivations for their actions. Thanks to such an ability, agents would become not only able to acquire novel information but also novel capabilities (i.e., procedural knowledge), possibly acquiring the (bits of) self-awareness required to provide *explanations* about their own courses of action.

In practice, the imagination abstraction leverages mechanisms laying at the intersection of different research areas, such as: (i) symbolic reasoning and automatic theorem proving (which are exploited to synthesize novel knowledge via inference), (ii) automatic planning and simulation (which is exploited to speculate over alternative courses of action), (iii) machine learning and data mining (which are exploited to induce new knowledge from experience), and (iv) biochemical coordination (which keeps imagination dynamic by continuously reorganizing it).

The rest of the paper is organized as follows. Section 2 briefly presents the current background technologies and their state-of-the-art supporting our notion of imagination. Section 3 introduces, defines, and discusses the concept of imagination and our practical view. Section 4 elicits the challenges related to our definition of imagination and the related research areas involved. Finally, Section 5 concludes the paper.

## 2 State of the Art

The investigation of mechanisms for agents’ imagination roots from cross-disciplinary components. In particular, this section provides a brief background on (i) human imagination mechanisms; (ii) cognitive agent architectures; (iii) imagination mechanisms including inference, data-driven learning, biochemical coordination, & simulation; and (iv) computational creativity.

### 2.1 Imagination in Humans

In the late 80s and early 90s, constructivist [6, 7] and developmental [8] approaches inspired many advancements in AI. In particular, virtual agents have been equipped with “inherent” learning mechanisms, allowing them *to make sense* of their environment and exploit its affordances<sup>4</sup> [10, 6]. Such an approach

<sup>4</sup> The term, originally coined by Gibson, refers to what the environment offers an individual [9]

has been inspired by constructivism [11] and developmental learning, promoted by the Swiss psychologist Jean Piaget (1930s), and adapted by Drescher into a bottom-up developmental approach to interact with the surrounding environment named **Schema** mechanism [6].

In the context of autonomous agents, such **Schemas** are pieces of knowledge processed by the agents to comprehend and react to their environments. The developmental process of **Schemas** is characterized by:

**Assimilation:** describes how humans or agents perceive and adapt to new information, fitting them into existing cognitive schemas.

**Accommodation:** restructures existing **Schemas** to handle novel information.

Intuitively, the cognitive growth of an intelligent system would imply the evolution/extension of both very specific knowledge and overall dynamics coordinating the several learning-related aspects. In the last decade, constructivist approaches have been used for smart environments [12, 13] and transport systems [14]. The studies contributing to these aspects provide contributions highly specialized in simplistic and very structured domains [15]. Nevertheless, such individual approaches can be hand-crafted together into architectures relying on traditional component-based software development methodologies. However, although the single components can evolve singularly, the reconciling system is constrained by the the hand-crafted interconnection mechanisms. Consequently, such lack of flexibility precludes architecture-level evolution (i.e., autonomous architectural adaptation and growth of the systems) and learning [7].

The lack of generalization characterizing current solutions impedes the application of intelligent/learning systems in general-purpose scenarios, being incapable of applying themselves autonomously to arbitrary problems. Therefore, imagination cannot be a cross-system functionality.

## 2.2 Cognitive Agents

The philosopher Michael Bratman formalized human practical reasoning in the beliefs-desires-intentions (BDI) model as a way to explaining future-directed intention [16]. Successively, it became a model to program intelligent agents, which made its first appearance in the Rational Agency project at the Stanford Research Institute in the mid 1980s [17]. Such a model is characterized by:

**Beliefs:** a set of facts and rules representing an agent’s epistemic memory, possibly containing its knowledge about the world, itself, and other agents.

**Desires:** a set of goals the agent is willing to achieve, test, or maintain.

**Intentions:** a set of tasks the agent is currently carrying on.

**Plans:** a set of *recipes* representing the agent’s procedural memory, encoding the procedural know-how about tasks.

Any cognitive feature of a BDI agent may vary during its lifetime. For instance, novel beliefs appear in the agents’ minds whenever they receive novel perceptions from their sensors, while stale beliefs simultaneously disappear. Similarly, novel beliefs may arise while agents interact among each other – or with

humans – or as they chose to memorize some information they have deduced via reasoning. The occurrence of relevant events may provoke the desire pool’s update (i.e., acquiring new goals to be achieved/tested/maintained and/or discard some goals). Agents’ desires eventually lead to spawning novel intentions (activities to achieve/test/maintain goals the agent is committed to). While carrying on an intention, agents may select one or more plans among those supporting the corresponding desire’s accomplishment. Plans may involve the execution of one or more actions – possibly affecting the world via actuators – or the accomplishment of further sub-goals, which may, in turn, require the execution of further plans as part of the same intention.

In the scope of this paper, it is worth highlighting that the BDI model allows agents to exhibit more complex behavior than purely reactive models, unbound to the computational overhead of other cognitive architectures [17]. Furthermore, being rooted in folk-psychology, it has been outlined as an excellent candidate to represent everyday explanations [18, 19] (since it is considered as the attribution of human behavior using “everyday” terms such as beliefs, desires, intentions, emotions, and personality traits [20, 21]).

The BDI model has also been identified as the most used/suitable architecture to generate explanations for goal-driven agents/robots [19, 22, 23]. The trend of attributing to the BDI model the suitability for XAI applications is reinforced by user studies supporting the human tendency to attribute a State of Mind (SoM) to robots and agents. In such a context, the lack of communication or misalignment due to lack of transparency can result in ill-formed SoM [24]. To avoid such a risk and the consequent drop of trust in the system, BDI agents are envisioned to employ folk-psychology to explain their SoM [19, 25].

### 2.3 Mechanisms for Imagination

In our view, the agents’ imagination process must rely on mechanisms employing different techniques, such as inference, data-driven learning, biochemistry-inspired coordination, and simulation.

**Logic Inference** In computational logic [26], inference is the process of rigorously drawing conclusions out of premises. The existing inference procedures depend on the given logic formalism at hand, and they may serve different purposes (depending on their nature). Overall, there are three main sorts of inference: deductive, inductive, and abductive.

**Deductive inference** dictates under which conditions conclusions can be drawn out of some axioms, i.e., rules and facts considered true. In other words, deduction elicits the knowledge which is possibly implicit into the axioms.

**Inductive inference** aims at estimating rules out of a number of positive (and negative) examples of facts satisfying (or violating) the rule. In other words, induction attempts to generalize principles by distilling patterns from generic observations/contingencies.

**Abductive inference:** aims at hypothesizing which premises could provoke some evidences, given a number of rules describing how causes provoke effects. In other words, abduction attempts to speculate on the possible causes of some phenomenon (e.g., finding the most straightforward and most likely explanation for an observation), given that the general rules governing that phenomenon are known.

Logic programming (LP) technologies (e.g., Prolog) enable users – and potentially agents – to encode their knowledge into logic facts and rules, which may then be queried via logic solvers [27]. Accordingly, by endowing agents with adequate LP technologies, they can autonomously exploit inference when required [27].

**Learning from Data** Data-driven AI falls into the context of the so-called machine learning (ML). Learning from data is commonly the activity performed via supervised or unsupervised learning and comprises a broad set of methods and tools such as reinforcement learning, classification, regression, time series forecasting, pattern recognition, generative models [28]. Supervised learning leverages on the existence of many input/expected-output examples and consists of looking for the best function mapping the available inputs into the corresponding expected outputs. In a sense, supervised learning is very similar to logic induction, except that it does not assume knowledge to be encoded via logic clauses, and it is better suited for *learning* from numeric data. Unsupervised learning aims at finding similarities and patterns possibly buried into numerical data without any expected outcome at hand. Thus, pieces of information are extracted from data through some optimality criterion.

By exploiting the wide availability of task-specific techniques and algorithms in ML, agents may be equipped with the capability of managing different sorts of data to serve disparate purposes [29]. For instance, by wrapping neural networks, agents may gain image and speech recognition capabilities, as well as the capability of analyzing and forecasting time-related measurements.

Finally, it is worth mentioning another relevant perspective intersecting ML and MAS: learning autonomously, continuously, and adaptively to increment skills and knowledge (a.k.a, lifelong ML or continuous learning—CL henceforth). In the context of ML, it entails updating the prediction models periodically with novel tasks and data distributions, still being able to (re)use and retain knowledge and skills over time. CL is beneficial when data or tasks’ availability varies over time (i.e., no longer or not yet available), and it is imperative to consider prior knowledge [30, 31].

**Biochemical Coordination** Within the scope of self-organizing MAS, biochemical coordination is the study of interaction among agents mediated by biochemistry-inspired patterns. There, information is modelled as *molecules*, i.e., chunks of data characterized by a *concentration* value denoting their relevance [32]. Such molecules may *diffuse* among different locations (e.g., to represent information exchanges), *aggregate* with each others (e.g., to represent

more complex data structures), and *evaporate*, (e.g., reduce their concentration as the carried information loses relevance). The concentration and nature of such molecules determine the dynamics of the systems relying on the biochemical metaphor. A number of coordination rules are commonly in place, affecting (and being affected by) the concentration of molecules within a given context, and governing information diffusion, aggregation, evaporation, or generation.

Due to their nature, such sorts of systems are inherently stochastic and fuzzy, and therefore ideal to realize resilient, robust, and self-organizing applications. In this context, pieces of information are not solely true or false, but rather more or less *concentrated*. Therefore, inconsistencies and contradictory data may simultaneously co-exist with no harm, as long as consistent truths eventually *emerge* by becoming significantly more concentrated. The combination of such features determines biochemical coordination mechanisms eligible to support imagination, as it may spawn several (possibly inconsistent) ideas, properly balancing evaporation, diffusion, and aggregation to retain only the most useful ones.

**Multi-agent Based Simulation** Simulation is one of the most employed techniques to identify/reach potentially useful outcomes. Agent simulation technology has been outlined as an efficient platform helping to understand autonomous behavior and decision-making [33]. An agent-based simulation (ABS) model is a set of interacting intelligent entities that reflect, within an artificial environment, the relationships in the real world [33]. Thus, ABS is typically used for helping decision-makers cope with complex and changing environment in the domains such as UAVs (c.f. [34] and the references therein), IoT, and CPS [35, 36], and to model and optimize robot behavior [37].

It is worth noticing that most of the works in the ABS literature focus on inter-agent relationships and their interaction with the environment [17]. Conversely, this paper aims at tackling the intra-agent perspective, where agents should be capable of simulating multiple states of themselves and their actions within their own “mind”. This internal simulation process is analogous to human “mental simulation” where humans rely on the ability to construct mental models to imagine what will happen or what could be [38–40]. Such capability has helped humans in physical reasoning [41, 42], spatial reasoning, and counterfactual reasoning [43].

Similarly, agents can mimic this “mental” modeling and analyze the assumed outcomes of its own actions, identify and possibly exploring arguably reasonable paths leading to potentially creative scenarios even in robustly novel situations [40]. Such explorations might lead to totally unforeseen solutions, which without a simulation based on a trial and error approach would not have been discovered/investigate. Hence, agents may reflect upon a set of simulations representing themselves in different/alternative scenarios.

## 2.4 Computational Creativity

In its broadest scope, creativity is defined as the ability to generate new forms and artifacts autonomously [44]. In the literature, creativity is classified as either

*biological* (the ability to generate new cells, organs, organisms, or species [44]) or *psychological* (the ability to generate new ideas and artifacts). Researches in AI have been pushing to extend the notion of creativity to virtual systems [45]. For example, a recent study, inspired from enactive AI [46, 47], investigates how computational creativity can be adopted by autonomous agents [48]. Despite this progress, most of the works in this domain are either carried out at the conceptual level or solely rely on data-driven mechanisms (e.g., generative adversarial networks, a.k.a. GANs) to generate “creative” contents (e.g., music [49] or pictures [50]).

In contrast with these works (primarily ML-centered), we envision agents questioning their beliefs, knowledge, and goals continuously. In particular, agents should combine classic planning, reinforcement learning, and in-mind simulation about their future actions to simulate and possibly provide explanations about their courses of actions.

### 3 Imagination in Cognitive Agents

Overall, BDI agents’ dynamics are moved by intentions and directed by desires. Equipped with *sensors* and *actuators*, they can respectively *perceive* and *affect* the world they live into. However, an agent’s admissible pool of desires and plans is defined/constrained by human developers. Indeed, developers tend to dictate agents’ initial desires and plans to keep their dynamics predictable and controllable. However, this prevents the full exploitation of agents’ autonomy, adaptability, and, ultimately, intelligence.

Arguably, to let agents access a higher degree of intelligence, they should be endowed with the capability of spontaneously generating new desires, acquiring novel beliefs, and learning novel plans. Briefly speaking, we consider imagination as the activity devoted to supporting such capability. Thus, we define imagination as an agent’s intention aimed at *maintaining* its *innate* desires of being *creative*, *curious*, and *effective*. More precisely, in our framework, agents are assumed to be endowed with (at least) one *maintenance* desire since their creation. Such desire pushes them to (attempt to) be creative, curious, and effective whenever they can. To be creative, an agent should keep looking for novel information, as well as novel ways to do what it needs to do (i.e., it must keep trying to enrich its belief and plan bases). To be curious, an agent should keep exploring the world and search for novel stimuli or just doing things to learn something new (i.e., it must keep trying to enrich its desires). To be effective and prove the way it deliberates and acts, an agent should keep improving its epistemic and procedural knowledge (i.e., improve its belief and plan bases).

To accomplish such an innate desire, BDI agents must spawn an intention that will be part of them for their whole lifetime. The basic functioning of this intention is relatively straightforward: *to keep revising the agents’ beliefs, desires, and plans to generate novel epistemic/procedural knowledge or improve the current one*. We call this intention “imagination”.

To accomplish its purpose, the imagination intention may leverage and combine several basic mechanisms coming from different branches of AI. Imagination can exploit mechanisms deriving from the classes of activities listed below (independently from technical details). For example,

**knowledge acquisition** is the process of converting raw data (i.e., percepts or beliefs) into general and reusable knowledge (e.g., in the form of logic rules or sub-symbolic predictors)

**knowledge synthesis** is the process of inferring or distilling novel knowledge out of pre-existing ones

**speculation** is the process of exploring alternative truths, situations, or courses of actions based on previous experiences

**knowledge revision** is the process of criticizing the pre-existing knowledge, possibly evicting stale or wrong information

The remainder of this section analyzes how mechanisms from the many branches of AI may be exploited to support such activities. Figure 1 provides a summarizing characterization.

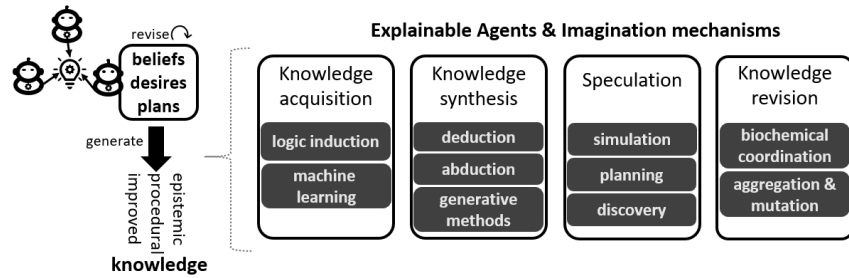


Fig. 1. Imagination in cognitive agents: AI mechanisms.

### 3.1 Acquiring Knowledge via Learning and Induction

BDI agents acquire information either by perceiving the environment or by communicating with other agents. In both cases, information comes in the form of raw data and can be *stored* either symbolically or sub-symbolically. Each datum represents a particular event from the external world. While the single event may be potentially useless *per se*, the frequent occurrence of similar events may generate value in the long run. Indeed, agents – similarly to humans – may distill valuable *knowledge* out of statistically relevant anterior experience (i.e., data).

Differently from data, however, knowledge is an *aggregated* and *reusable* form of information. It must be reusable because that is what makes it valuable enough to memorize it. It must be aggregated because agents’ cognitive resources (such as computational power and memory) are inherently limited, in practice, and

such limitations affect how and to what extent data can be actually reused. However, the way data is actually aggregated to make it reusable depends remarkably on its nature.

In case data is symbolically represented, it is interpreted as logic *facts* and stored into symbolic knowledge bases that the agent may efficiently update and query. When this is the case, logic *induction* may be used to distill *rules* out of facts. While facts are contingent, rules are synthetic, and they may be reused in several similar situations. Furthermore, symbolic rules are human-intelligible. Thus, they can be used by agents as a basis to construct explanations for their reasoning or behavior.

Conversely, when data is represented sub-symbolically, machine learning can be exploited to draw knowledge out of it. When this is the case, data is interpreted as tensors of numbers used to *train* a predictor (e.g., a neural network). This usually makes knowledge both aggregated and reusable, despite not directly intelligible (and explainable) for humans. Accordingly, induction can be exploited by agents willing or requiring to manipulate symbolic information, either because they need to take discrete decisions or because they care about the intelligibility of their decisional process. Conversely, machine learning can be exploited by agents needing to manipulate sub-symbolic information—possibly because they need to take fuzzy decisions, and can tolerate errors to a certain extent.

### 3.2 Synthesizing Knowledge via Deduction, Abduction, and Generative Methods

The external world is not the only source of valuable knowledge. Indeed, intelligent agents should also be able to *synthesize* novel knowledge out of what they already know. The way they do so, however, may vary depending on the nature of the knowledge at hand. For instance, when knowledge is represented in symbolic form, *deductive* or *abductive* reasoning procedures may be exploited to infer novel information out of it. Conversely, when knowledge is sub-symbolically represented, *generative* methods may be exploited instead.

In particular, deductive reasoning may be exploited to make *implicit* knowledge explicit. In fact, deduction can derive specific facts out of rules. In other words, it is dual w.r.t. induction. Thanks to deduction, agents may for instance select useful knowledge for the contingent situation they are immersed into, out of general rules. Similarly, abductive reasoning may be exploited by agents willing to draw *hypotheses* about the causes which lead to a particular situation. In other words, abduction let agents synthesize likely facts which justify the facts they already know to be true, according to the rules they know to hold in a particular context. Accordingly, abduction is one of the mechanisms supporting speculative thinking.

Conversely, generative methods – such as GAN – may be exploited by agents needing to produce human-comprehensible representations (e.g. audio, video, etc.) of categories they already know how to recognize and manipulate—such as faces [51], shapes [52], handwriting [53], speech, etc. These representations might

serve as key enablers for explainable synthesized knowledge. In turn, generative methods may be key enablers for *(i)* computational creativity, as they let agents produce original representations; *(ii)* counterfactual thinking, as they support the generation of *variants* of any given concept; and *(iii)* effective human-machine interactions, as they let agents enrich their interactions with humans with randomly-generated examples or analogies.

### 3.3 Speculating via Simulation and Planning

Mentally simulating scenarios is a fundamental human capability [39]. Once people have enough information about the characterization and dynamics of the surrounding world, simulating the effects of their actions becomes a common practice (to a certain extent). Often, it is only after having *mentally* simulated the most likely outcomes of their course of actions that an individual chooses how to act. Then, by comparing the *actual* outcomes with the expected ones, humans may learn how to improve their behavior w.r.t. their goals. Furthermore, even when a direct experience is lacking, simulating the possible effects of a given action is still better than acting randomly.

In the AI literature, *planning* is the activity performed by agents willing to deliberate what to do in a particular context. Planning algorithms commonly leverage on rich descriptions of *(i)* the environment, *(ii)* agents' actions, *(iii)* their effects, and *(iv)* some description of the target goal the agent is willing to achieve. Through such descriptions, planning algorithms (attempt to) compute viable workflows of actions that should lead agents towards the target goal. However, even when only a few agents and small deterministic environments are involved, planning is computationally costly. Therefore, when complex non-deterministic environments are in place (where several agents interact in non-trivial ways), planning may quickly become unfeasible.

Scientific researchers often tackle the complexity of systems by *simulating* simplified parametric models executed multiple times, with randomly generated parameters. Doing so allows drawing statistical conclusions based on the data generated by such *in-silico* experiments. Accordingly, agents may follow a similar approach, in their minds, to decide what to do or what to expect. In other words, agents may leverage simulation to realize *speculative* thinking.

MAS have been exploited for the purpose of simulation since their very beginning—cf. ABS in Section 2.3. However, currently, most ABS research efforts are devoted to exploiting MAS *in* simulations rather than the opposite. Conversely, the idea of letting an agent simulate itself and its environment is quite new—other than very challenging. According to such a perspective, we envision equipping each agent with an ABS sub-system capable of simulating an entire MAS. Using *inner* simulations, agents could then try out different actions and sequence of actions that are otherwise too costly or even dangerous to try in the real environment, other than retrying the same scenario over and over again with different rules or parameters. Similarly to humans, agents may then exploit such capability to autonomously discover plans, rules, or even policies for situations that they have either already experienced, or not.

### 3.4 Adaptively Revising Knowledge via Biochemical Coordination

Simulating the world to discover novel plans, rules, or – more generally – background knowledge may eventually lead to the creation of chunks of *potentially* useless information. By doing so, efficiency issues in both information storage and retrieval may arise. However, a more conservative strategy may prevent agents from discovering novel and *potentially* useful information. Accordingly, some general strategies should be exploited to allow each agent to decide what knowledge to retain and discard dynamically.

Here we welcome the idea that no “one size fits all” solution exists to select knowledge based on expected utility. Indeed, any predefined strategy may be affected by the biases of who designed it or be tailored to a particular scenario while being sub-optimal in other ones. For this reason, we argue that an *adaptive* strategy based on a biochemical metaphor would be preferable.

From such a perspective, we assume the many mechanisms proposed so far – inference, machine learning, planning, simulation – to produce information in the form of molecules. Agents’ minds can then be conceived as containers of molecules of different sorts—e.g., beliefs, neurons, plans, etc. Such molecules’ concentration may increase over time as the corresponding information may be produced multiple times or on a per-usage basis. For instance, the same rule may be induced from different data in different instants or be frequently used. Similarly, different runs of a simulation may lead to the frequent exploitation of similar courses of action. This, in turn, may lead an agent to increase the concentration of one or more plans. At the same time, we assume the information is subject to *evaporation* on a uniform basis. In other words, all sorts of information evaporate at the same pace. As a global effect, only relevant information would be able to survive, in the long run – where by relevant we mean either frequently generated or frequently used, without requiring to *a-priori* define what is actually relevant.

Finally, aggregation and mutation mechanisms may fit the picture by supporting random and periodic modifications or combinations of pre-existing information—e.g., by merging similar plans/rules into more articulated ones, or by slightly altering some parameter of a wrapped neural network to change its behavior. Such random modifications may then be adaptively confirmed or discarded, depending on whether they result to be relevant or not for the agent. In the former case, their concentration would increase, whereas in the latter case, it would decrease.

## 4 Open Challenges

As we discussed in the previous sections, imagination within the context of agent intentions entails the provision of continuous knowledge acquisition, synthesis, revision and exploration. While we identified both the existing approaches (cf. Section 2) as well as the specific research areas (cf. Section 3) to investigate, we acknowledge the existence of key challenges to tackle in the process. These are summarized as follows:

- C1) *Knowledge heterogeneity.* Acquisition of knowledge is essential to feed a creative process, whether it is intending to combine information to derive new insights, for synthesizing explainable knowledge, for confronting different views, or even for exploration of uncharted territory. Nevertheless, agents will be exposed to the challenge of extreme variety among the different knowledge sources that they run across. Solving semantic and data-representation heterogeneity issues arising from this diversity will be a necessary step. An example would be using knowledge graph matching and fusion techniques [54]. Moreover, given the need for integrating symbolic and sub-symbolic sources, tools and techniques for a coherent integration between both will have to be studied [55].
- C2) *Goal generation.* A fundamental step in a creative cycle is establishing clear goals, even if these may be updated in the future. While a goal may define an overall scope for the development of creative activities, in some cases, the goals may not be entirely known *a priori*. In such conditions, goal generation [56] must be part of the creative process that needs to be incorporated into the agent model [57].
- C3) *Knowledge alignment.* Even when knowledge heterogeneity has been addressed, it will still be required to align different understandings of observed phenomena relevant to the creative plan's scope. For example, if knowledge sources' provenance is dissimilar, simply aligning terminologies and concepts is not enough [58]. At this point, it is crucial to study models that allow handling contradictions, assumptions, explainable outcomes, and interpretations as part of the knowledge alignment task [59, 60].
- C4) *Information uncertainty.* Creative agents must take into account not only the potential inaccuracy of their information sources but also the eventual uncertainty of their own artificial imagination. The exploration and navigation over radically new ideas and approaches entail high risk, meaning that oftentimes they may lead to dead-ends. Agents may need to incorporate risk management strategies [61] allowing them not to constrain themselves only to *safe* knowledge but leave enough space for behavioral models that adapt to different levels of uncertainty. This also applies to uncertainty in XAI outcomes and their consequences on inter-agent agreements.
- C5) *Reasoning complexity.* The generation of new knowledge may require reasoning over potentially large and/or complex knowledge graphs [62]. Depending on the complexity of these graphs' underlying logics, reasoning tasks may become increasingly expensive in terms of computation. Moreover, the agents' autonomous nature will necessitate further exploration of decentralized reasoning techniques, including partial knowledge and probabilistic approaches. An additional challenge to tackle is the combination of explicable results of data-driven AI predictions. Multi-agent speculative reasoning may need to be combined with machine learning outcomes to address this challenge.
- C6) *Hypotheses evaluation.* Agents will be able to propose hypotheses that may need to be validated or refuted [63]. This ability should be accompanied by a robust framework for managing assumptions, claims, justifications, explanations, and proofs [64, 65]. As explained in the previous point, reasoning and

sub-symbolic outcomes have to be evaluated with respect to the hypotheses. Agents may eventually have different or plainly contradictory points of view, for which reconciling mechanisms may need to put in place. While in some cases competitive approaches may be preferred (e.g., working towards the same goal but under different imaginative hypotheses and assumptions), in others, it might be necessary to align and establish a cooperation scheme.

- C7) *Explicable knowledge revision.* When the results of explicable machine learning and, in general, generated sub-symbolic knowledge are produced, agents need to navigate through them and understand their implication over existing information [66]. This may lead to invalidating previous beliefs or to changing the uncertain status or certain facts. The challenge of explaining these decisions, and providing justification of the imaginative paths taken by a community of agents, shall be addressed to understand the path that leads to creative activities. The provenance of knowledge and the changes may lead to even reconsidering information that was deemed false or invalid in a previous iteration.
- C8) *Exploration.* Agent imagination requires substantial space for the exploration of new knowledge and experimentation through novel approaches. Although exploratory agents have been studied in the past [67], it remains a challenge to establish a formal framework for discovery in large knowledge spaces. Approaches like link traversal of knowledge graphs may serve as a starting point, although they may need to be extended to a cooperative scenario where different agents run exploratory tasks under coordination mechanisms.
- C9) *Accountability.* Imaginative processes in multi-agent systems entail the exploration and creation of new knowledge, as well as the validation of previous and new ideas. The consequences of these actions may lead to decisions and actions for which there should be clear responsibilities. In that context, the provenance information emanating from exploratory processes and knowledge revision decisions will need to be associated with trust mechanisms allowing to ensure proper attribution to an agent or a person embodied by an agent. Furthermore, accountability [68] in terms of ethical and even legal terms should be studied, not only from a purely technical perspective (e.g., accountability networks, knowledge graph ledgers) but also from a psycho-social point of view (i.e., human-agent accountability).

## 5 Conclusions

This paper provided a ground for discussing the meaning of imagination in the setting of cognitive agents, selected possible tools and approaches, and elicited the envisioned contextual challenges. In particular, the investigated research are *(i)* reasoning and automatic theorem proving, *(ii)* automatic planning and simulation, *(iii)* machine learning and data mining, and *(iv)* biochemical coordination. Finally, the intuitions proposed directions collapsed in the definition of challenges in the areas of knowledge heterogeneity, goal generation/definition, knowledge alignment, information uncertainty, reasoning complexity, hypothesis evaluation, explicable knowledge revision, exploration, and accountability.

To address these challenges, in the future we plan to explore a number of practical research directions aimed at creating the technological playground for supporting our notion of imagination. For instance, the problem of letting agent programming technologies support several logics and inference procedures, is far from being solved [69]. A similar statement holds for the simulation of large-scale MAS composed by cognitive agents. For this reason, our first efforts shall be devoted to *(i)* the development (resp. extension) of novel (resp. existing) agent programming framework to support inductive, and abductive reasoning – for instance, via the 2P-KT technology [70] –, *(ii)* the development of simulation frameworks for cognitive agents, supporting virtualization of both space and time – for instance via the Alchemist simulator [71] –, *(iii)* blending (either existing or novel) agent programming frameworks and mainstream ML frameworks—such as TensorFlow, PyTorch, etc. Conversely, concerning the design of and biochemical coordination at the single agent level, we argue that further research is needed. Along this path, our first step will consist of a formalization, aimed at further clarifying possibilities and challenges.

## Acknowledgements

This work has been partially supported by the CHIST-ERA grant CHIST-ERA-19-XAI-005, and by *(i)* the Swiss National Science Foundation (G.A. 20CH21\_195530), *(ii)* the Italian Ministry for Universities and Research, *(iii)* the Luxembourg National Research Fund (G.A. INTER/CHIST/19/14589586), *(iv)* the Scientific and Research Council of Turkey (TÜBİTAK, G.A. 120N680).

## References

1. Dietrich Stout. Stone toolmaking and the evolution of human culture and cognition. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 366(1567):1050–1059, 2011.
2. Stephen Cave. The problem with intelligence: Its value-laden history and the future of ai. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 29–35, 2020.
3. Vasant Honavar. Symbolic artificial intelligence and numeric artificial neural networks: towards a resolution of the dichotomy. In *Computational Architectures integrating Neural and Symbolic Processes*, pages 351–388. Springer, 1995.
4. Giovanni Ciatto, Michael I. Schumacher, Andrea Omicini, and Davide Calvaresi. Agent-based explanations in AI: Towards an abstract framework. In Davide Calvaresi, Amro Najjar, Michael Winikoff, and Kary Främling, editors, *Explainable, Transparent Autonomous Agents and Multi-Agent Systems*, volume 12175 of *Lecture Notes in Computer Science*, pages 3–20. Springer, Cham, 2020. Second International Workshop, EXTRAAMAS 2020, Auckland, New Zealand, May 9–13, 2020, Revised Selected Papers.
5. Rafael H Bordini and Jomi F Hübner. Bdi agent programming in agentspeak using jason. In *International workshop on computational logic in multi-agent systems*, pages 143–164. Springer, 2005.

6. Gary L Drescher. *Made-up minds: a constructivist approach to artificial intelligence*. MIT press, 1991.
7. Kristinn R Thórisson. From constructionist to constructivist ai. In *AAAI Fall Symposium: Biologically Inspired Cognitive Architectures*, 2009.
8. Douglas Blank, Deepak Kumar, Lisa Meeden, and James B Marshall. Bringing up robot: Fundamental mechanisms for creating a self-motivated, self-organizing architecture. *Cybernetics and Systems: An International Journal*, 36(2):125–150, 2005.
9. James Jerome Gibson. *The Senses Considered as Perceptual Systems*. Houghton Mifflin, 1966.
10. Rodney A Brooks. Intelligence without representation. *Artificial intelligence*, 47(1-3):139–159, 1991.
11. Jean Piaget. *La naissance de l’intelligence chez l’enfant*. paris: Delachaux et Niestlé, 1936.
12. Amro Najjar and Patrick Reignier. Constructivist ambient intelligent agent for smart environments. In *2013 IEEE International Conference on Pervasive Computing and Communications Workshops (PERCOM Workshops)*, pages 356–359. IEEE, 2013.
13. Amr Alzouhri Alyafi, Van-Bao Nguyen, Yann Laurillau, Patrick Reignier, Stephane Ploix, Gaëlle Calvary, Joëlle Coutaz, Monalisa Pal, and Jean-Philippe Guilbaud. From usable to incentive building energy management systems. *Modélisation et utilisation du contexte (Modeling and Using Context)*, 2(1):1–30, 2018.
14. Maxime Guériaux, Frédéric Armetta, Salima Hassas, Romain Billot, and Nour-Eddin El Faouzi. A constructivist approach for a self-adaptive decision-making system: application to road traffic control. In *2016 IEEE 28th International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 670–677. IEEE, 2016.
15. Olivier L Georgeon, Jonathan H Morgan, and Frank E Ritter. An algorithm for self-motivated hierarchical sequence learning. In *Proceedings of the International Conference on Cognitive Modeling. Philadelphia, PA. ICCM-164*, pages 73–78. Citeseer, 2010.
16. Michael Bratman et al. *Intention, plans, and practical reason*, volume 10. Harvard University Press Cambridge, MA, 1987.
17. Carole Adam and Benoit Gaudou. Bdi agents in social simulations: a survey. *The Knowledge Engineering Review*, 31(3):207–238, 2016.
18. Emma Norling. Folk psychology for human modelling: Extending the bdi paradigm. In *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems-Volume 1*, pages 202–209, 2004.
19. Joost Broekens, Maaike Harbers, Koen Hindriks, Karel Van Den Bosch, Catholijn Jonker, and John-Jules Meyer. Do you get it? user-evaluated explainable bdi agents. In *German Conference on Multiagent System Technologies*, pages 28–39. Springer, 2010.
20. Paul M Churchland. Folk psychology and the explanation of human behavior. *Philosophical Perspectives*, 3:225–241, 1989.
21. Bertram F Malle. How people explain behavior: A new theoretical framework. *Personality and social psychology review*, 3(1):23–48, 1999.
22. Mark A Neerincx, Jasper van der Waa, Frank Kaptein, and Jurriaan van Diggelen. Using perceptual and cognitive explanations for enhanced human-agent team performance. In *International Conference on Engineering Psychology and Cognitive Ergonomics*, pages 204–214. Springer, 2018.

23. Sule Anjomshoae, Amro Najjar, Davide Calvaresi, and Kary Främling. Explainable agents and robots: Results from a systematic literature review. In *18th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2019), Montreal, Canada, May 13–17, 2019*, pages 1078–1088. International Foundation for Autonomous Agents and Multiagent Systems, 2019.
24. Thomas Hellström and Suna Bensch. Understandable robots-what, why, and how. *Paladyn, Journal of Behavioral Robotics*, 9(1):110–123, 2018.
25. Frank Kaptein, Joost Broekens, Koen Hindriks, and Mark Neerincx. Personalised self-explanation by robots: The role of goals versus beliefs in robot-action explanation for children and adults. In *2017 26th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, pages 676–682. IEEE, 2017.
26. Lawrence C. Paulson. Computational logic: its origins and applications. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2210):20170872, 2018.
27. Roberta Calegari, Giovanni Ciatto, Enrico Denti, and Andrea Omicini. Logic-based technologies for intelligent systems: State of the art and perspectives. *Information*, 11(3):1–29, March 2020. Special Issue “10th Anniversary of Information—Emerging Research Challenges”.
28. Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Pearson, 2002.
29. Giovanni Ciatto, Roberta Calegari, Andrea Omicini, and Davide Calvaresi. Towards XMAS: eXplainability through Multi-Agent Systems. In Claudio Savaglio, Giancarlo Fortino, Giovanni Ciatto, and Andrea Omicini, editors, *AI&IoT 2019 – Artificial Intelligence and Internet of Things 2019*, volume 2502 of *CEUR Workshop Proceedings*, pages 40–53. Sun SITE Central Europe, RWTH Aachen University, November 2019.
30. Bing Liu. Lifelong machine learning: a paradigm for continuous learning. *Frontiers of Computer Science*, 11(3):359–361, 2017.
31. Yuwei Cui, Subutai Ahmad, and Jeff Hawkins. Continuous online sequence learning with an unsupervised neural network model. *Neural computation*, 28(11):2474–2504, 2016.
32. Jose Luis Fernandez-Marquez, Giovanna Di Marzo Serugendo, Sara Montagna, Mirko Viroli, and Josep Lluís Arcos. Description and composition of bio-inspired design patterns: a complete overview. *Natural Computing*, 12(1):43–67, 2013.
33. Michael J Wooldridge and Nicholas R Jennings. Intelligent agents: Theory and practice. *The knowledge engineering review*, 10(2):115–152, 1995.
34. Yazan Mualla, Wenshuai Bai, Stéphane Galland, and Christophe Nicolle. Comparison of agent-based simulation frameworks for unmanned aerial transportation applications. *Procedia computer science*, 130:791–796, 2018.
35. Davide Calvaresi, Mauro Marinoni, Arnon Sturm, Michael Schumacher, and Giorgio Buttazzo. The challenge of real-time multi-agent systems for enabling iot and cps. In *Proceedings of the international conference on web intelligence*, pages 356–364, 2017.
36. Davide Calvaresi, Giuseppe Albanese, Jean-Paul Calbimonte, and Michael Schumacher. Seamless: Simulation and analysis for multi-agent system in time-constrained environments. In *International Conference on Practical Applications of Agents and Multi-Agent Systems*, pages 392–397. Springer, 2020.
37. Daniel Urieli, Patrick MacAlpine, Shivaram Kalyanakrishnan, Yinon Bentor, and Peter Stone. On optimizing interdependent skills: a case study in simulated 3d humanoid robot soccer. In *AAMAS*, volume 11, page 769, 2011.

38. Philip N Johnson-Laird. Inference with mental models. *The Oxford handbook of thinking and reasoning*, pages 134–145, 2012.
39. Dedre Gentner and Albert L Stevens. Mental models lawrence erlbaum associates. *Hillsdale, New Jersey*, 1983.
40. Jessica B Hamrick. Analogues of mental simulation and imagination in deep learning. *Current Opinion in Behavioral Sciences*, 29:8–16, 2019.
41. Mary Hegarty. Mechanical reasoning by mental simulation. *Trends in cognitive sciences*, 8(6):280–285, 2004.
42. Peter W Battaglia, Jessica B Hamrick, and Joshua B Tenenbaum. Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45):18327–18332, 2013.
43. Paul L Harris. *The work of the imagination*. Blackwell Publishing, 2000.
44. Margaret A Boden et al. *The creative mind: Myths and mechanisms*. Psychology Press, 2004.
45. Margaret A Boden. Creativity and alife. *Artificial Life*, 21(3):354–365, 2015.
46. Tom Froese and Tom Ziemke. Enactive artificial intelligence: Investigating the systemic organization of life and mind. *Artificial Intelligence*, 173(3-4):466–500, 2009.
47. Pierre De Loor, Kristen Manac’h, and Jacques Tisseau. Enaction-based artificial intelligence: Toward co-evolution with humans in the loop. *Minds and Machines*, 19(3):319–343, 2009.
48. Christian Guckelsberger, Christoph Salge, and Simon Colton. Addressing the “why?” in computational creativity: A non-anthropocentric, minimal model of intentional creative agency. In *International Conference on Computational Creativity (ICCC 2017)*, pages 128–135. Goldsmiths, University of London, 2017.
49. Olof Mogren. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*, 2016.
50. Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. GauGAN: semantic image synthesis with spatially adaptive normalization. In *ACM SIGGRAPH 2019 Real-Time Live!*, pages 2:1–2:10, 2019.
51. Angelo Genovese, Vincenzo Piuri, and Fabio Scotti. Towards explainable face aging with generative adversarial networks. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 3806–3810. IEEE, 2019.
52. Carlo Biffi, Juan J Cerrolaza, Giacomo Tarroni, Wenjia Bai, Antonio De Marvao, Ozan Oktay, Christian Ledig, Loic Le Folgoc, Konstantinos Kamnitsas, Georgia Doumou, et al. Explainable anatomical shape analysis through deep hierarchical generative models. *IEEE transactions on medical imaging*, 39(6):2088–2099, 2020.
53. Yao Zhu, Saksham Suri, Pranav Kulkarni, Yueru Chen, Jiali Duan, and C-C Jay Kuo. An interpretable generative model for handwritten digits synthesis. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 1910–1914. IEEE, 2019.
54. Michael Azmy, Peng Shi, Jimmy Lin, and Ihab F Ilyas. Matching entities across different knowledge graphs with graph embeddings. *arXiv preprint arXiv:1903.06607*, 2019.
55. Artur S d’Avila Garcez, Krysia B Broda, and Dov M Gabbay. *Neural-symbolic learning systems: foundations and applications*. Springer Science & Business Media, 2012.
56. Jan Broersen, Mehdi Dastani, Joris Hulstijn, and Leendert van der Torre. Goal generation in the boid architecture. *Cognitive Science Quarterly*, 2(3-4):428–447, 2002.

57. Zhizhou Ren, Kefan Dong, Yuan Zhou, Qiang Liu, and Jian Peng. Exploration via hindsight goal generation. *arXiv preprint arXiv:1906.04279*, 2019.
58. Jérôme Euzenat. Interaction-based ontology alignment repair with expansion and relaxation. In *IJCAI 2017-26th International Joint Conference on Artificial Intelligence*, pages 185–191. AAAI Press, 2017.
59. Paula Chocron and Paolo Paretì. Vocabulary alignment for collaborative agents: a study with real-world multilingual how-to instructions. In *IJCAI*, pages 159–165, 2018.
60. Ernesto Jiménez-Ruiz, Terry R Payne, Alessandro Solimando, and Valentina Tamma. Limiting logical violations in ontology alignment through negotiation. In *Proceedings of the Fifteenth International Conference on Principles of Knowledge Representation and Reasoning*, pages 217–226, 2016.
61. Martin Lorenz, Jan D Gehrke, Hagen Langer, Ingo J Timm, and Joachim Hammer. Situation-aware risk management in autonomous agents. In *Proceedings of the 14th ACM international conference on Information and knowledge management*, pages 363–364, 2005.
62. Luigi Bellomarini, Eleonora Laurenza, Emanuel Sallinger, and Evgeny Sherkhonov. Reasoning under uncertainty in knowledge graphs. In *International Joint Conference on Rules and Reasoning*, pages 131–139. Springer, 2020.
63. Fatima B Seeme and David G Green. Pluralistic ignorance: Emergence and hypotheses testing in a multi-agent system. In *2016 International Joint Conference on Neural Networks (IJCNN)*, pages 5269–5274. IEEE, 2016.
64. Ioan Alfred Letia and Adrian Groza. Arguing with justifications between collaborating agents. In *International Workshop on Argumentation in Multi-Agent Systems*, pages 102–116. Springer, 2011.
65. Jonatan Olofsson, Gustaf Hendeby, Tom Rune Lauknes, and Tor Arne Johansen. Multi-agent informed path planning using the probability hypothesis density. *Autonomous Robots*, pages 1–13, 2020.
66. Roberto Confalonieri, Tillman Weyde, Tarek R Besold, and Fermín Moscoso del Prado Martín. Using ontologies to enhance human understandability of global post-hoc explanations of black-box models. *Artificial Intelligence*, page 103471, 2021.
67. Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *International Conference on Machine Learning*, pages 2778–2787. PMLR, 2017.
68. Matteo Baldoni, Cristina Baroglio, Olivier Boissier, Katherine Marie May, Roberto Micalizio, and Stefano Tedeschi. Accountability and responsibility in agent organizations. In *International Conference on Principles and Practice of Multi-Agent Systems*, pages 261–278. Springer, 2018.
69. Roberta Calegari, Giovanni Ciatto, Viviana Mascardi, and Andrea Omicini. Logic-based technologies for multi-agent systems: A systematic literature review. *Autonomous Agents and Multi-Agent Systems*, 35(1):1:1–1:67, 2021. Collection “Current Trends in Research on Software Agents and Agent-Based Software Development”.
70. Giovanni Ciatto, Roberta Calegari, Enrico Siboni, Enrico Denti, and Andrea Omicini. 2P-KT: logic programming with objects & functions in kotlin. In Roberta Calegari, Giovanni Ciatto, Enrico Denti, Andrea Omicini, and Giovanni Sartor, editors, *WOA 2020 – 21th Workshop “From Objects to Agents”*, volume 2706 of *CEUR Workshop Proceedings*, pages 219–236, Aachen, Germany, October 2020. Sun SITE Central Europe, RWTH Aachen University. 21st Workshop “From Objects to Agents” (WOA 2020), Bologna, Italy, 14–16 September 2020. Proceedings.

20 G. Ciatto, A. Najjar, J.P. Calbimonte, and D. Calvaresi

71. Danilo Pianini, Sara Montagna, and Mirko Viroli. Chemical-oriented simulation of computational systems with ALCHEMIST. *Journal of Simulation*, 2013.