# Multi-Agent Causal Reinforcement Learning

André Meyer-Vitali[a]

*Deutsches Forschungszentrum für Künstliche Intelligenz GmbH (DFKI), Saarbrücken, Germany*

Keywords: Software Engineering, Artificial Intelligence, Trust, Transparency, Robustness, Causality, Agency.

Abstract: It has become clear that mere correlations extracted from data through statistical processes are insufficient to give insight into the causal relationships inherent in them. Causal models support the necessary understanding of these relationships to make transparent and robust decisions. In a distributed setting, the causal models that are shared between agents improve their coordination and collaboration. They learn individually and from each other to optimise a system's behaviour. We propose a combination of causal models and multi-agent reinforcement learning to create reliable and trustworthy AI systems. This combination strengthens the modelling and reasoning of agents that communicate and collaborate using shared causal insights. A comprehensive method for applying and integrating these aspects is being developed.

## 1 INTRODUCTION

The development of reliable and trustworthy AI systems requires new methods and metrics for their creation and verification. However, not everything needs to be developed from scratch, because the experience of decades of systems and software engineering, as well as research in the foundations of AI systems, provides ample opportunities to reuse and adapt existing methods. We outline some principles of AI Engineering in section 2. We dive deeper into two of the principles, namely agency in section 3 and causality in section 4. Consequently, in section 5, the combination of causality and multi-agent reinforcement learning (MARL) is investigated. A use case concerning urban mobility and energy consumption is presented in section 6. Finally, some conclusions and an outlook are provided in section 7.

## 2 PRINCIPLES OF TRUSTED AI

The design and development of complex systems has a long-standing tradition in software architecture and engineering (Gamma et al., 1994; Booch et al., 2005). This experience can be used to build better AI systems, including the current wave of data-driven and hybrid systems, that are reliable and worthy of trust.

**Trust** is defined as the willingness of a trustor to be vulnerable to the actions of an actor (trustee) that she cannot directly control. It requires the existence of uncertainty and risks. The level of trust increases with the trustor's perceived benevolence, competence and integrity of the trustee (Lewis and Weigert, 1985; Rousseau et al., 1998; Mayer et al., 1995; Jacovi et al., 2021; Lewis and Marsh, 2022; de Brito Duarte et al., 2023). Trust calibration is the process of adjusting and aligning actual and perceived trustworthiness, where actual trustworthiness is the degree to which a system or actor complies with its required or promised expectations and performance (Okamura and Yamada, 2020; de Visser et al., 2020; Visser et al., 2023; Chi and Malle, 2023)).

In order to achieve trustworthiness and trust, an AI system should include a number of characteristics, such as those listed by the High-Level Expert Group (HLEG) of the EU. Their Assessment List for Trustworthy Artificial Intelligence, ALTAI (Directorate-General for Communications Networks, Content and Technology (European Commission), 2020)) lists:

1. Human Agency and Oversight;
2. Technical Robustness and Safety;
3. Privacy and Data Governance;
4. Transparency;
5. Diversity, Non-discrimination and Fairness;
6. Societal and Environmental Well-being;
7. Accountability.

[a] https://orcid.org/0000-0002-5242-1443

An analysis of AI engineering methods for trust (Meyer-Vitali and Mulder, 2024a; Meyer-Vitali and Mulder, 2024b) reveals the following four main principles.

**Models & Explanations.** Reliable predictions and decisions about system behaviour for insightful and plausible explanations and simulations with generalised models from knowledge and training.

**Causality & Grounding.** Identification and predictions of cause-effect relationships for informed predictions and actions, as well as anchoring of meaning in real-world context and phenomena.

**Modularity & Compositionality.** Design of complex systems broken down into comprehensible and manageable parts (functions and features), reliably composed in system architectures.

**Human Agency & Oversight.** Overview, final decision, and human responsibility for the actions of AI systems, also when delegating tasks to autonomous agents in hybrid collaborative teams.

An important aspect to consider is that we should not aim to reinvent the wheel, time and time again. Instead, we should augment what is already established knowledge and build on top of existing models and methods (Schreiber et al., 1999; Tiddi et al., 2023). Furthermore, it is not sufficient to determine statistical patterns and correlations in data. What machine "learning" achieves is just that – which is not learning as a desire to understand. Beyond "learning", we should aim for understanding concepts and relationships. AI can become a tool for scientific discovery if developers and systems understand what they learned and manage to convert experimental insights into hypotheses and theories that serve as stepping stones for further experiments. This is the scientific method that was successful for centuries and which should not be given up for fancy trompe l'oeils (Goldstein and Goldstein, 1978; Gower, 1996; Nola and Sankey, 2007; Griffin et al., 2024; Jamieson et al., 2024).

## 3 AGENCY

Agency is the level of control that an entity has over itself and its behaviour (van der Vecht et al., 2007). Agents are "systems that can decide for themselves what they need to do in order to satisfy their design objectives" (Wooldridge in (Weiss, 2000)). Thus, agents reason and take deliberate decisions on their behalf and initiative (Wooldridge and Jennings, 1995; Wooldridge, 2009). Agents are fundamental building blocks of AI systems (Russell and Norvig, 2020;

OECD, 2022), which were investigated for a long period and are currently of major interest, again, for the future of AI (Larsen et al., 2024).

In technical terms, we are interested in agents as autonomous and communicative entities (software or robots) that learn, reason, plan and act to achieve some goals. Each agent can communicate with its environment (physical or virtual) and with other agents, including humans. Due to their autonomy, agents can process data and knowledge efficiently in a distributed fashion, which also benefits scalability and sovereignty.

Agents don't come alone. In multi-agent systems (MAS) (Weiss, 2000) agents take roles and interact with each other to achieve their individual or collective goals. They communicate using communicative or speech acts (Searle, 1969)according to a variety of interaction patterns and protocols, such as negotiations, auctions or queries (Poslad and Charlton, 2001).

The benefits of multi-agent systems for building reliable AI systems include

- modularisation and compositionality by defining specific roles and protocols for communicating among agents;

- separation of concerns, where agents have their individual roles, values and duties;

- agreements are reached by negotiations, such that outcomes can be explained; and

- distribution of logic and processing for efficiency, scalability and sovereignty.

## 4 CAUSALITY

Causality refers to our ability to understand and predict how things are connected through cause and effect – specifically, what causes what, and why (Pearl and Mackenzie, 2018). When artificial intelligence systems can grasp these causal connections, they become capable of making better predictions and solving complex problems. This represents a crucial shift in focus: rather than simply identifying correlations (when things happen together), we need to understand causation (when one thing directly leads to another).

This transition from correlation-based to causation-based thinking is becoming increasingly crucial, particularly as we seek to understand the reasoning behind AI predictions and decisions (Pearl et al., 2016). Understanding the true causes behind outcomes, rather than just their associations, is essential for explaining why AI systems make specific choices.

Causal inference considers how and when causal conclusions can be drawn from data. Complex systems of interacting variables can be described with a Structural Causal Model (SCM) (Pearl, 2009; Peters et al., 2017; Nogueira et al., 2022). Such a model describes the causal mechanisms and assumptions present in an arbitrary system. The relationships between (endogenous) variables can be visualised in a directed acyclic graph (DAG), wherein nodes represent variables and directed paths represent causal influences between variables. These influences do not need to be deterministic, but can be probabilistic, include external factors (exogenous variables) and are valid for linear as well as non-linear functions, discrete as well as continuous variables.

A **structural causal model (SCM)** is a triple (Pearl, 2009)

$$M = \langle U, V, F \rangle,$$

where:

- $U$ is a set of background variables, (also called *exogenous*), that are determined by factors outside the model;

- $V$ is a set $\{V_1, V_2, \cdots, V_n\}$ of variables, called *endogenous*, that are determined by variables in the model;

- $F$ is a set of functions $\{f_1, f_2, \cdots, f_n\}$ such that each $f_i$ is a mapping from $U_i \cup PA_i{}^1$ to $V_i$: $v_i = f_i(pa_i, u_i), i = 1, \cdots, n$.

A causal model $M$ can be visualised as a directed acyclic graph, $G(M)$, in which each node corresponds to a variable and the directed edges point from members of $PA_i$ and $U_i$ toward $V_i$. This graph is the causal diagram associated with $M$. From a causal diagram, it is easy to read dependencies between variables by examining the three different types of connections – chains, forks, and colliders – using $d$-separation[2]. Graphs also help to identify confounders: variables that influence both the treatment (cause) and the outcome (effect). Confounders, as common causes, create a spurious association between treatment and outcome.

The **Law of Conditional Independence** (*d*-separation)

$$(X \, sep \, Y | Z)_{G(M)} \Rightarrow (X \perp\!\!\!\perp Y | Z)_{P(v)}$$

describes the consequence of separation in the model as independence in the distribution.

Interventions (or experiments) correspond to forcing a variable $X$ to take on value $x$, thereby removing dependencies on parent variables to examine changes

in one variable (representing a state, action or event) and whether they cause changes in another, in order to distinguish between correlated and causal relationships in data.

An **intervention $I$** in an SCM $M$ entails changing some set of structural assignments in $M$ with a new set of structural assignments. Assume the replacement is on $X_k$ given by assignment $X_k = \tilde{f}(\tilde{PA}_k, \tilde{U}_k)$, where $\tilde{PA}_k$ are the parents in the new graph.

Counterfactuals refer to alternative choices that could have been made in the past and the corresponding effects that they might have caused. Therefore, they allow for exploring possibilities to find alternative outcomes according to a causal model, allowing to change policies accordingly in the future.

The **Law of Counterfactuals**

$$Y_x(u) = Y_{M_x}(u)$$

states that a model $M$ generates and evaluates all counterfactuals. $M_x$ is a modified structural submodel of $M$, where the equation for $X$ is replaced by $X = x$. This allows for generating answers to an enormous number of hypothetical questions of the type "What would $Y$ be had $X$ been $x$?" (Pearl et al., 2016).

An SCM can be interpreted as a probability distribution $P$ with density $p$ over variables $X$ in the causal system. According to the *Ladder of Causation* (Pearl and Mackenzie, 2018), three classes of reasoning exist.

1. **Seeing** (associations) encapsulates statistical reasoning.

$$P(y|x)$$

2. **Doing** (interventions) contains randomised control trials (RCTs) and RL methods (cf. section 5).

$$P(y|do(x), z)$$

3. **Imagining** (counterfactuals) allows for reasoning about outcomes of alternative choices.

$$P(y_x|x', y')$$

Consequently, this formalism allows for explicit reasoning about each action an agent *does*, *can*, and *could* have taken. Interestingly, causality provides means to understand the story behind data, i.e., a causal model can be seen as a data generator that reflects the causal relationships. Therefore, it can explain those relationships and is seen as a necessary requirement for explainability (Carloni et al., 2023; Rawal et al., 2023; Ganguly et al., 2023). An SCM can be used for planning, such that causes are determined to achieve desired effects (Meyer-Vitali and Mulder, 2023).

---

[1]$PA_i$ denotes the parents of $V_i$.

[2]$d$ stands for directional.

# 5 MULTI-AGENT CAUSAL REINFORCEMENT LEARNING

Reinforcement Learning (RL) (Sutton and Barto, 2018) is a method for agents to learn how to map situations to actions from interacting with an environment. Agents act and receive feedback on the quality or performance of their actions (rewards), which they use for updating their strategy or policy.

Typically, the strategy is defined as a Markov Decision Process (MDP). Markov decision processes include three aspects: sensation, action, and goal. MDPs consist of a finite set of states $S$, a finite set of actions $A$, a reward function $R : S \times A \times S \to \mathbb{R}$, a state transition probability function $T : S \times A \times S \to [0,1]$ and an initial state distribution.

More realistically, not all relevant information can be sensed or observed from the environment and the set of states and actions is infinite. Therefore, Partially Observable MDPs (POMDP) are often used (Kaelbling et al., 1998). POMDPs do not sense the states directly, but receive observations which depend (in a probabilistic or deterministic way) on the state of the environment. Agents will need to take into account the history of past observations to infer the possible current state of the environment.

In a distributed setting, Multi-Agent Reinforcement Learning (MARL) concerns multiple agents learning concurrently or collaboratively in the same environment (Buşoniu et al., 2010; Nowé et al., 2012; Zhang et al., 2021; Gronauer and Diepold, 2022; Albrecht et al., 2024). For example, in Partially Observable Stochastic Games (POSG) all agents have the same reward function, joint observations and a joint action set. They are also known as "Decentralized POMDP" (Dec-POMDP) and are used in the area of multi-agent planning (Oliehoek, 2012; Oliehoek and Amato, 2016). However, in a generalised approach, agents do not need to obey to the same reward function for distributed learning and control.

Finally, agents, causality and reinforcement learning come together. The reader may appreciate that the similarity between interventions and counterfactuals in causal models with agents' actions and rewards in reinforcement learning provides a useful synergy. An approach at this combination was already presented by (Maes et al., 2007; Grimbly et al., 2021; Jiao et al., 2024), where multi-agent causal models are introduced. Agents share an environment and have access to private and public variables of interest.

A Multi-Agent Causal Model (MACM) consists of n agents, each of which contains a semi-Markovian model $M_i$:

$$M_i = \langle V_{M_i}, G_{M_i}, P(V_{M_i}), K_{M_i} \rangle, i \in \{1, \cdots, n\}$$

- $V_{M_i}$ is the subset of variables that agent $a_i$ can access.

- $G_{M_i}$ is the causal graph over variables $V_{M_i}$.

- $P(V_{M_i})$ is the joint probability distribution over $V_{M_i}$.

- $K_{M_i}$ stores the intersections $V_{M_i,M_j}$ with other agents $a_j$, $\{V_{M_i} \cap V_{M_j}\}$, assuming that the agents agree on the structure and distribution of their intersections.

MACMs are useful as shared causal models in teams of agents. Each agent reasons according to its individual POMDP with endogenous and exogenous variables and corresponding observations of the environmental state. Agents share histories of observations and rewards (state-action trajectories), but take individual coordinated actions based on their individual rewards. An MACM can be seen as a generalisation of Dec-POMDPs.

Multi-Agent Causal Reinforcement Learning (MACRL) allows for causal inference in a dynamic learning environment with a multitude of agents (Casini and Manzo, 2016; Pina et al., 2023a; Pina et al., 2023b; Richens and Everitt, 2024). With further elaboration of these basic ideas, MACRL could evolve into several design patterns for neuro-symbolic AI systems (van Bekkum et al., 2021), where neuro-causal agents combine explicit prior knowledge in the form of (hypothetical) causal models with discovery (Pearl, 2019) and learning for continuous self-improvement.

# 6 SCENARIO: URBAN AGENTS

Urban life has many peculiar characteristics (WBGU – German Advisory Council on Global Change, 2016; Angelidou et al., 2022; Oliveira et al., 2020; Popelka et al., 2023; Petrikovičová et al., 2022; Hashem et al., 2023). Many different streams of information, activities and resources are intertwined and have conflicting requirements and characteristics, such as energy, mobility, food, water, waste, healthcare, commerce, and many more (Nevejan et al., 2018). They all share the same limited spatial and temporal constraints. The urban context requires and involves many interactions at high speed with lots of people, high density and diversity of the population, land use, displacements (geographical, social and occupational mobility), as well as fluent and heterogeneous communities.

The interactions in an urban environment are diverse, complex and conflicting. Many interests of hybrid actors are related and depend on each other. In the urban context, an overall goal for sustainable use
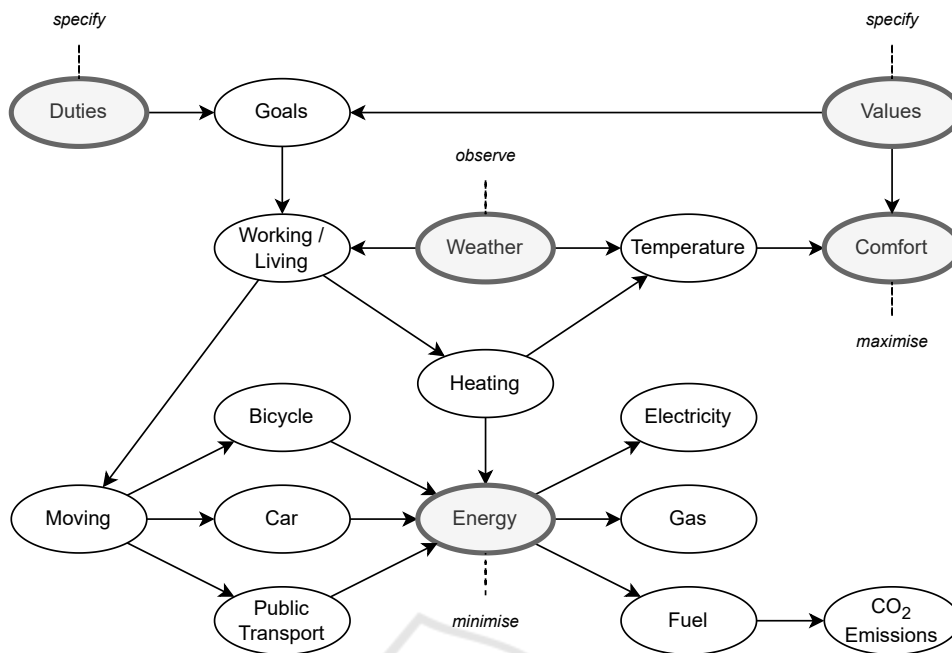
Figure 1: A causal model for living and working in the urban context.

of resources could be the reduction of energy consumption.

Some causal relationships in an urban context, focusing on a combination of energy consumption and mobility, are shown in figure 1. In the urban causal model, *values*, *duties* and the *weather* are exogenous variables. They need to be accepted as they are and are independent. *Goals*, *working/living* and *moving* are causes that determine the behaviour of agents. *Energy* and *comfort* are the effects that we want to achieve and to optimise.

The consumption of various types of energy is affected by the need and desire to move about the city and to heat buildings at home and at work (and for leisure, shopping, etc.). Values and duties are the main sources that drive urban behaviour and external factors, such as the weather, influence decision-making. This causal model explains the relationships among several important behavioural aspects, but it is not deterministic. Individual behaviour is influenced by exogenous variables and cooperative behaviour results in complex interactions.

Each of the variables is modelled using a quantitative metric. The transfer functions are to be determined using causal machine learning, based on the behaviour of actors.

A shared goal can be seen and modelled as an effect, that is caused by one or more interventions (actions or events). Consequently, in order to decide and plan which actions to take, it is necessary to understand which actions or events cause the intended effects. For example, your goal can be to arrive at a destination at a given time (work, home, leisure, etc.). By reasoning back which actions are required to get you there, piece by piece, a connected causal path can be constructed to determine the departure time and modes of traffic along the route. Due to shared intentions and causal models, humans and agents can mutually trust each other regarding their actions and outcomes.

# 7 CONCLUSIONS AND OUTLOOK

After looking at agency, causality and reinforcement learning separately, the combination in Multi-Agent Causal Reinforcement Learning (MACRL) provides a path towards more robust, transparent and explainable AI systems, such that they become more trustworthy (concerning some of the characteristics of the ALTAI). Model-based software engineering for such AI systems promotes reliability and trust, because the actions and interventions that agents take are clearly understandable. By investigating various ways of implementing MACRL, a more robust software engineering method may emerge. This paper gives some insights and directions for such a method; more research is needed and should be applied in real-world use cases.

# REFERENCES

Albrecht, S. V., Christianos, F., and Schäfer, L. (2024). *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press.

Angelidou, M., Politis, C., Panori, A., Bakratsas, T., and Fellnhofer, K. (2022). Emerging smart city, transport and energy trends in urban settings: Results of a pan-European foresight exercise with 120 experts. *Technological Forecasting and Social Change*, 183:121915.

Booch, G., Rumbaugh, J., and Jacobson, I. (2005). *Unified Modeling Language User Guide, The (2nd Edition) (Addison-Wesley Object Technology Series)*. Addison-Wesley Professional, 2nd edition.

Buşoniu, L., Babuška, R., and De Schutter, B. (2010). Multi-agent Reinforcement Learning: An Overview. In Srinivasan, D. and Jain, L. C., editors, *Innovations in Multi-Agent Systems and Applications - 1*, pages 183–221. Springer, Berlin, Heidelberg.

Carloni, G., Berti, A., and Colantonio, S. (2023). The role of causality in explainable artificial intelligence. arXiv:2309.09901 [cs].

Casini, L. and Manzo, G. (2016). *Agent-based models and causality: a methodological appraisal*. Linköping University Electronic Press.

Chi, V. B. and Malle, B. F. (2023). Calibrated Human-Robot Teaching: What People Do When Teaching Norms to Robots*. In *2023 32nd IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, pages 1308–1314. ISSN: 1944-9437.

de Brito Duarte, R., Correia, F., Arriaga, P., and Paiva, A. (2023). AI Trust: Can Explainable AI Enhance Warranted Trust? *Human Behavior and Emerging Technologies*, 2023:e4637678. Publisher: Hindawi.

de Visser, E. J., Peeters, M. M. M., Jung, M. F., Kohn, S., Shaw, T. H., Pak, R., and Neerincx, M. A. (2020). Towards a Theory of Longitudinal Trust Calibration in Human–Robot Teams. *International Journal of Social Robotics*, 12(2):459–478.

Directorate-General for Communications Networks, Content and Technology (European Commission) (2020). *The Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self assessment*. Publications Office of the European Union.

Gamma, E., Helm, R., Johnson, R., Vlissides, J., and Booch, G. (1994). *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison-Wesley Professional, Reading, Mass, 1st edition edition.

Ganguly, N., Fazlija, D., Badar, M., Fisichella, M., Sikdar, S., Schrader, J., Wallat, J., Rudra, K., Koubarakis, M., Patro, G. K., Amri, W. Z. E., and Nejdl, W. (2023). A Review of the Role of Causality in Developing Trustworthy AI Systems. arXiv:2302.06975 [cs].

Goldstein, M. and Goldstein, I. F. (1978). *How We Know: An Exploration of the Scientific Process*. Westview Press.

Gower, B. (1996). *Scientific Method: A Historical and Philosophical Introduction*. Routledge, London.

Griffin, C., Wallace, D., Mateos-Garcia, J., Schieve, H., and Kohli, P. (2024). A new golden age of discovery. Technical report, DeepMind.

Grimbly, S. J., Shock, J., and Pretorius, A. (2021). Causal Multi-Agent Reinforcement Learning: Review and Open Problems. arXiv:2111.06721 [cs].

Gronauer, S. and Diepold, K. (2022). Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, 55(2):895–943.

Hashem, I. A. T., Usmani, R. S. A., Almutairi, M. S., Ibrahim, A. O., Zakari, A., Alotaibi, F., Alhashmi, S. M., and Chiroma, H. (2023). Urban Computing for Sustainable Smart Cities: Recent Advances, Taxonomy, and Open Research Challenges. *Sustainability*, 15(5):3916. Number: 5 Publisher: Multidisciplinary Digital Publishing Institute.

Jacovi, A., Marasović, A., Miller, T., and Goldberg, Y. (2021). Formalizing Trust in Artificial Intelligence: Prerequisites, Causes and Goals of Human Trust in AI. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pages 624–635, New York, NY, USA. Association for Computing Machinery.

Jamieson, K. H., Kearney, W., and Mazza, A.-M., editors (2024). *Realizing the Promise and Minimizing the Perils of AI for Science and the Scientific Community*. University of Pennsylvania Press.

Jiao, L., Wang, Y., Liu, X., Li, L., Liu, F., Ma, W., Guo, Y., Chen, P., Yang, S., and Hou, B. (2024). Causal Inference Meets Deep Learning: A Comprehensive Survey. *Research*, 7:0467. Publisher: American Association for the Advancement of Science.

Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. (1998). Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134.

Larsen, B., Li, C., Teeuwen, S., Denti, O., DePerro, J., and Raili, E. (2024). Navigating the AI Frontier: A Primer on the Evolution and Impact of AI Agents. Technical report, World Economic Forum.

Lewis, J. D. and Weigert, A. (1985). Trust as a Social Reality. *Social Forces*, 63(4):967–985.

Lewis, P. R. and Marsh, S. (2022). What is it like to trust a rock? A functionalist perspective on trust and trustworthiness in artificial intelligence. *Cognitive Systems Research*, 72:33–49.

Maes, S., Meganck, S., and Manderick, B. (2007). Inference in multi-agent causal models. *International Journal of Approximate Reasoning*, 46(2):274–299.

Mayer, R. C., Davis, J. H., and Schoorman, F. D. (1995). An Integrative Model of Organizational Trust. *The Academy of Management Review*, 20(3):709–734. Publisher: Academy of Management.

Meyer-Vitali, A. and Mulder, W. (2023). Causing Intended Effects in Collaborative Decision-Making. In Murukannaiah, P. K. and Hirzle, T., editors, *Proceedings of the Workshops at the Second International Conference on Hybrid Human-Artificial Intelligence*, volume 3456 of *CEUR Workshop Proceedings*, pages 137–144, Munich, Germany. CEUR. ISSN: 1613-0073.

Meyer-Vitali, A. and Mulder, W. (2024a). Engineering Principles for Building Trusted Human-AI Systems. In Arai, K., editor, *Intelligent Systems and Applications*, pages 468–485, Cham. Springer Nature Switzerland.

Meyer-Vitali, A. and Mulder, W. (2024b). Human-AI Engineering for Adults. In *HHAI 2024: Hybrid Human AI Systems for the Social Good*, pages 228–240. IOS Press.

Nevejan, C., Sefkatli, P., and Cunningham, S. W. (2018). *City rhythm: logbook of an exploration*. Delft University of Technology, Multi Actor Systems, Amsterdam, first edition edition. OCLC: 1312647218.

Nogueira, A. R., Pugnana, A., Ruggieri, S., Pedreschi, D., and Gama, J. (2022). Methods and tools for causal discovery and causal inference. *WIREs Data Mining and Knowledge Discovery*, 12(2):e1449. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/widm.1449.

Nola, R. and Sankey, H. (2007). *Theories of Scientific Method: an Introduction*. Routledge, London.

Nowé, A., Vrancx, P., and De Hauwere, Y.-M. (2012). Game Theory and Multi-agent Reinforcement Learning. In Wiering, M. and van Otterlo, M., editors, *Reinforcement Learning: State-of-the-Art*, pages 441–470. Springer, Berlin, Heidelberg.

OECD (2022). OECD Framework for the Classification of AI systems. Technical report, OECD, Paris.

Okamura, K. and Yamada, S. (2020). Adaptive trust calibration for human-AI collaboration. *PLOS ONE*, 15(2):1–20. Publisher: Public Library of Science.

Oliehoek, F. A. (2012). Decentralized POMDPs. In Wiering, M. and van Otterlo, M., editors, *Reinforcement Learning: State-of-the-Art*, pages 471–503. Springer, Berlin, Heidelberg.

Oliehoek, F. A. and Amato, C. (2016). *A Concise Introduction to Decentralized POMDPs*. SpringerBriefs in Intelligent Systems. Springer International Publishing, Cham.

Oliveira, G. M., Vidal, D. G., and Ferraz, M. P. (2020). Urban Lifestyles and Consumption Patterns. In Leal Filho, W., Marisa Azul, A., Brandli, L., Gökçin Özuyar, P., and Wall, T., editors, *Sustainable Cities and Communities*, Encyclopedia of the UN Sustainable Development Goals, pages 851–860. Springer International Publishing, Cham.

Pearl, J. (2009). *Causality: Models, Reasoning and Inference*. Cambridge University Press, Cambridge, U.K. ; New York, 2nd edition edition. https://bayes.cs.ucla.edu/BOOK-2K/.

Pearl, J. (2019). The seven tools of causal inference, with reflections on machine learning. *Communications of the ACM*, 62(3):54–60.

Pearl, J., Glymour, M., and Jewell, N. P. (2016). *Causal Inference in Statistics: A Primer*. John Wiley & Sons. Google-Books-ID: I0V2CwAAQBAJ.

Pearl, J. and Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect*. Basic Books, New York, 1st edition edition.

Peters, J., Janzing, D., and Schölkopf, B. (2017). *Elements of Causal Inference: Foundations and Learning Algorithms*. The MIT Press.

Petrikovičová, L., Kurilenko, V., Akimjak, A., Akimjaková, B., Majda, P., Ďatelinka, A., Biryukova, Y., Hlad, L., Kondrla, P., Maryanovich, D., Ippolitova, L., Roubalová, M., and Petrikovič, J. (2022). Is the Size of the City Important for the Quality of Urban Life? Comparison of a Small and a Large City. *Sustainability*, 14(23):15589. Number: 23 Publisher: Multidisciplinary Digital Publishing Institute.

Pina, R., De Silva, V., and Artaud, C. (2023a). Causality Detection for Efficient Multi-Agent Reinforcement Learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, AAMAS '23, pages 2824–2826, Richland, SC. International Foundation for Autonomous Agents and Multiagent Systems.

Pina, R., De Silva, V., and Artaud, C. (2023b). Learning Independently from Causality in Multi-Agent Environments. arXiv:2311.02741 [cs].

Popelka, S., Narvaez Zertuche, L., and Beroche, H. (2023). Urban AI Guide. Technical report, Zenodo.

Poslad, S. and Charlton, P. (2001). Standardizing Agent Interoperability: The FIPA Approach. In Luck, M., Mařík, V., Štěpánková, O., and Trappl, R., editors, *Multi-Agent Systems and Applications: 9th ECCAI Advanced Course, ACAI 2001 and Agent Link's 3rd European Agent Systems Summer School, EASSS 2001 Prague, Czech Republic, July 2–13, 2001 Selected Tutorial Papers*, pages 98–117. Springer, Berlin, Heidelberg.

Rawal, A., Raglin, A., Sadler, B. M., and Rawat, D. B. (2023). Explainability and causality for robust, fair, and trustworthy artificial reasoning. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications V*, volume 12538, pages 493–500. SPIE.

Richens, J. and Everitt, T. (2024). Robust agents learn causal world models. arXiv:2402.10877 [cs].

Rousseau, D. M., Sitkin, S. B., Burt, R. S., and Camerer, C. (1998). Introduction to Special Topic Forum: Not so Different after All: A Cross-Discipline View of Trust. *The Academy of Management Review*, 23(3):393–404. Publisher: Academy of Management.

Russell, S. and Norvig, P. (2020). *Artificial Intelligence: A Modern Approach*. Pearson, Hoboken, NJ, 4th edition.

Schreiber, G., Akkermans, H., Anjewierden, A., de Hoog, R., Shadbolt, N. R., Van de Velde, W., and Wielinga, B. J. (1999). *Knowledge Engineering and Management: The CommonKADS Methodology*. The MIT Press.

Searle, J. R. (1969). *Speech Acts: An Essay in the Philosophy of Language*. Cambridge University Press, Cambridge.

Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA.

Tiddi, I., De Boer, V., Schlobach, S., and Meyer-Vitali, A. (2023). Knowledge Engineering for Hybrid Intelligence. In *Proceedings of the 12th Knowledge Capture Conference 2023*, K-CAP '23, pages 75–82, New

York, NY, USA. Association for Computing Machinery.

van Bekkum, M., de Boer, M., van Harmelen, F., Meyer-Vitali, A., and Teije, A. t. (2021). Modular design patterns for hybrid learning and reasoning systems. *Applied Intelligence*, 51(9):6528–6546.

van der Vecht, B., Meyer, A. P., Neef, M., Dignum, F., and Meyer, J.-J. C. (2007). Influence-Based Autonomy Levels in Agent Decision-Making. In Noriega, P., Vázquez-Salceda, J., Boella, G., Boissier, O., Dignum, V., Fornara, N., and Matson, E., editors, *Coordination, Organizations, Institutions, and Norms in Agent Systems II*, Lecture Notes in Computer Science, pages 322–337, Berlin, Heidelberg. Springer.

Visser, E. d., Momen, A., Walliser, J., Kohn, S., Shaw, T., and Tossell, C. (2023). Mutually Adaptive Trust Calibration in Human-AI Teams. In Murukannaiah, P. K. and Hirzle, T., editors, *Proceedings of the Workshops at the Second International Conference on Hybrid Human-Artificial Intelligence*, volume 3456 of *CEUR Workshop Proceedings*, pages 188–193, Munich, Germany. CEUR. ISSN: 1613-0073.

WBGU – German Advisory Council on Global Change (2016). Humanity on the move: Unlocking the transformative power of cities. Technical report, WBGU, Berlin. Frauke Kraas, Claus Leggewie, Peter Lemke, Ellen Matthies, Dirk Messner, Nebojsa Nakicenovic, Hans Joachim Schellnhuber, Sabine Schlacke, Uwe Schneidewind.

Weiss, G. (2000). *Multiagent Systems: A Modern Approach to Distributed Artificial Intelligence*. MIT Press, Cambridge, MA, USA, 1st edition.

Wooldridge, M. (2009). *An Introduction to MultiAgent Systems*. John Wiley & Sons, 2nd edition.

Wooldridge, M. and Jennings, N. R. (1995). Intelligent agents: theory and practice. *The Knowledge Engineering Review*, 10(2):115–152.

Zhang, K., Yang, Z., and Başar, T. (2021). Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms. arXiv:1911.10635 [cs, stat].