

# Chapter 1

## The Social Concept of Trust as Enabler for Robustness in Open Self-Organising Systems

Gerrit Anders, Hella Seebach, Jan-Philipp Steghöfer, Wolfgang Reif,  
Elisabeth André, Jörg Hähner, Christian Müller-Schloer, and Theo Ungerer

**Abstract** The participants in open self-organising systems, including users and autonomous agents, operate in a highly uncertain environment in which the agents' benevolence cannot be assumed. One way to address this challenge is to use computational trust. By extending the notion of trust as a qualifier of relationships between agents and incorporating trust into the agents' decisions, they can cope with uncertainties stemming from unintentional as well as intentional misbehaviour. As a consequence, the system's robustness and efficiency increases. In this context, we show how an extended notion of trust can be used in the formation of system structures, algorithmically to mitigate uncertainties in task and resource allocation, and as a sanctioning and incentive mechanism. Beyond that, we outline how the

---

G. Anders (✉) • H. Seebach

Institute for Software and Systems Engineering, University of Augsburg, Augsburg, Germany  
e-mail: [anders@isse.de](mailto:anders@isse.de); [seebach@isse.de](mailto:seebach@isse.de)

J.-P. Steghöfer

Department of Computer Science and Engineering, Chalmers University of Technology |  
University of Gothenburg, Gothenburg, Sweden  
e-mail: [jan-philipp.steghofer@cse.gu.se](mailto:jan-philipp.steghofer@cse.gu.se)

W. Reif

Institute for Software and Systems Engineering, University of Augsburg, Augsburg, Germany  
e-mail: [reif@isse.de](mailto:reif@isse.de)

E. André

Human-Centered Multimedia, University of Augsburg, Augsburg, Germany  
e-mail: [elisabeth.andre@informatik.uni-augsburg.de](mailto:elisabeth.andre@informatik.uni-augsburg.de)

J. Hähner

Organic Computing Group, University of Augsburg, Augsburg, Germany  
e-mail: [jorg.hahner@informatik.uni-augsburg.de](mailto:jorg.hahner@informatik.uni-augsburg.de)

C. Müller-Schloer

Institute of Systems Engineering, University of Hannover, Hannover, Germany  
e-mail: [cms@sra.uni-hannover.de](mailto:cms@sra.uni-hannover.de)

T. Ungerer

Systems and Networking Group, University of Augsburg, Augsburg, Germany  
e-mail: [theo.ungerer@informatik.uni-augsburg.de](mailto:theo.ungerer@informatik.uni-augsburg.de)

users' trust in a self-organising system can be increased, which is decisive for the acceptance of these systems.

**Keywords** Computational trust • Uncertainty • Self-organisation • Open MAS • Robustness

## 1.1 Trust as a Measure of Uncertainty in Open Self-Organising Systems

In open self-organising systems, different participants, such as autonomous agents, human users, and other systems, work together with a strong influence of the environment. These participants communicate and cooperate at runtime in unforeseeable ways and do not always follow the intent of the system designers. They can pursue different goals, and it cannot be assumed that they are intrinsically motivated to contribute towards a common system goal [1, 2]. Beyond that, a participant's behaviour can vary over time. As there is also limited knowledge about and control over the behaviour of the participants in the system, only weak assumptions about them can be made – in particular, we have to abandon assumptions of benevolence of the autonomous agents. The system participants therefore have to deal with both unintentional as well as intentional misbehaviour of others. This situation is aggravated by additional factors that increase uncertainties as they influence the system in unpredictable ways. These factors comprise the environment, other systems the agents interact with, or the users. Another form of openness often regarded in multi-agent systems (MAS) research is present when agents can arbitrarily enter and leave the system [3]. Especially in safety- or mission-critical domains, such as manufacturing or power management, these challenges have to be taken very seriously.

In this chapter, we argue that trust – as a measure of uncertainty – is a key concept for achieving robustness and efficiency in open self-organising systems. The classic notion of *computational trust* in the MAS community is focused on the credibility of agents, i.e. the degree to which they fulfil their commitments. This view stems mainly from psychological and sociological research [4] and boils down to the selection of interaction partners in order to maximise the utility of individual interactions. Economic [5, 6] and computer science [7, 8] literature characterise trust as instrumental to manage *expectations* about others. In computer science, the term “computational trust” is used to stress that the trust in a system or a system's part, such as an agent, is assessed by means of a well-defined metric. Since both (a part of) the system or a human being can act in the role of the trustor, we can differentiate between system-to-system and user-to-system trust. Often, a strong connection between trust and risk is emphasised [9] since interactions that incur a high risk for the participating agents require a high expectation of the others' willingness to contribute in a beneficial manner. An empirically justified expectation reduces the *uncertainty* about the behaviour of another agent [10]. In computing systems, this is often captured by a numerical *trust value* [11].

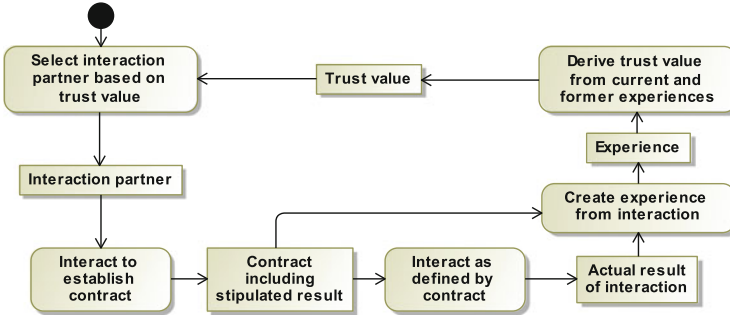
For these reasons, trust is an essential constituent of ensembles of cooperating agents, be they human or technical systems. Game-theoretical considerations show that trust can help to avoid getting trapped in the tragedy of the commons. Kantert et al. [12] provide such lines of thoughts in the context of Desktop Grid Computing. In general, trust induces a probability distribution over types of interaction partners of different trustworthiness in a Bayesian game. In this setting, agents have to choose their actions given probabilistic knowledge about each other's trustworthiness.

As mentioned above, we claim that trust is a key concept for achieving robustness. In this chapter, we define robustness in two dimensions. The first dimension of robustness addresses a system's ability to resist internal or external disturbances. Such disturbances result from (un)intentional misbehaving agents, for instance. A system exhibiting this type of robustness promises to remain in acceptable states and thus to maintain its functionality despite detrimental influences. The second dimension of robustness considers a system's ability to return into an acceptable state after a disturbance occurred that caused the system to leave the acceptance space. This type of robustness characterises a system's ability to restore its functionality. Consequently, the magnitude of disturbances the system can cope with (first dimension) and the duration of the deviation from acceptable states (second dimension) can be used to quantify the robustness. Both dimensions of robustness quantify the system's ability to fulfil its tasks. In contrast to a mere passive resistance, self-organising systems can actively increase their robustness by means of reactive or proactive measures. In open systems, these measures can be based on participants' trustworthiness, which allows the system to anticipate different sources of uncertainties.

In this chapter, we give an overview of the uses of computational trust (see Sect. 1.3) to deal with uncertainties arising in open self-organising systems. We show that these uses extend the classical use of selecting interaction partners and are based on the same life-cycle describing how trust values evolve over time (see Sect. 1.2). In detail, we demonstrate how trust models can be used to inform self-organisation processes (see Sect. 1.3.1); to optimise for critical or likely situations in uncertain environments (see Sect. 1.3.2); to sanction or incentivise agents in normative systems (see Sect. 1.3.3); and to represent the social relationships of the system's users (see Sect. 1.3.4). Section 1.4 concludes the chapter by emphasising that trust proves to be very useful to increase robustness and efficiency in open self-organising systems.

## 1.2 Computational Trust

Trust is usually measured as a numerical value, often normalised to values between 0 and 1. In [13], an agent's trust value is either very high or very low if the agent is either always expected to behave beneficially or never; if the value is between these extremes, the agent behaves in an unpredictable fashion and thus interactions with it are afflicted with a high uncertainty. Such a simple representation of trust is used



**Fig. 1.1** The life-cycle of trust values derived from experiences (adapted from [OCT3])

in many trust models (for an overview, see, e.g. [14]). However, numerous other interpretations and representations of trust exist. Anders et al. [OCT1], for instance, regard a trust value as an expected deviation from a prediction or promise. The lower an agent's trust value, the higher the expected deviation from its predictions or promises. A supplementary value, called *predictability*, quantifies the variance in the agent's behaviour and is used to indicate the certainty that the expected deviation actually occurs. Other representations based on more complex data structures (e.g. trust-based scenarios [OCT2] or elaborate reputation systems [15]) are able to capture further properties, such as time-dependent behaviour in the sense that an agent's behaviour depends on the time of day or that its behaviour depends on those it showed in previous time steps. Before discussing the general properties of trust, we illustrate the life-cycle of trust values which can be transferred to most of the other representations of trust.

*The Life-Cycle of Trust Values.* There is a general way of thinking about the origin of trust values that is independent of the way they are used (see Fig. 1.1). Two or more parties commit to a (potentially implicit) *contract* [16] that defines an *interaction* (possibly composed of several distinct steps) as well as its *stipulated result*. The *actual result* of the interaction can be compared to what was stipulated in the contract, thus yielding an *experience* for each party [17]. Ultimately, an agent uses its experiences and a trust metric to derive a *trust value* for each of its interaction partners. The trust values, in turn, inform future interactions.

Falcone et al. [18] criticised that many trust models are void of semantics of how the generated trust values have to be interpreted. It is, e.g. often not defined what a trust value of, say, 0.5 actually expresses or which trust value should be assigned to a new agent (the problem of *initial trust*, see, e.g. [19]). If a trust model has precise semantics, meaning a clearly defined way to interpret generated trust values, such an abstracting quantification can still be valid, though.

*Properties of Trust.* The life-cycle shows why trust values are *subjective*. As each agent makes its own experiences with others, it forms a personal opinion (i.e. a trust value) based on these unique experiences. Thus, the experiences of two agents with the same partner can vary tremendously. Additionally, agents can use different metrics to assess trust values and apply different requirements to the behaviour of others, thus implementing different trust models. The same arguments can be used to argue against *transitivity* of trust [20]. An exception are recommendations as a form of *indirect trust* or *reputation* (see discussion below) that have to be based on a mutual understanding of the valuation of an agent's behaviour.

Further, it is crucial to consider the *context* in which interactions occur. The context includes, e.g. the roles the agents play in the interaction, its contract, or environmental circumstances. Comparing experiences to each other in different contexts is difficult: You might trust your doctor to fix you, but not necessarily to fix your car. Falcone et al. [18] relate to context when they mention the “competence belief” an agent has about another. Competence is specific to a certain goal that the trusting agent believes the potential partner is capable to pursue. Agents that are deemed competent for one goal can be incompetent for another. Other authors use, e.g. “circumstance” [21] or “domain of interaction” [22] to denote context.

A trust value can also be supplemented by a measure of *confidence* [OCT4] or *certainty* [23, 24] that indicates the degree of certainty that a trust value describes the actual observable behaviour of an agent. Such an additional value can be based on several criteria, such as how many experiences were used for the calculation of the trust value, how old these experiences are, or how much the experiences differed. It is also possible to take the social relationships between the agents into account [25] or to distinguish short-term and long-term behaviour in order to identify changing behaviour. As with trust values themselves, the initialisation of confidence can be problematic. In human interactions, different trust dispositions are common where people approach newcomers differently and are willing to put more initial trust in them than others [26]. The experiences made by these trusting individuals can then be used by others to judge newcomers. Such a mechanism is especially useful during the exploratory phase after the start of a system [19].

*Reputation.* In open self-organising systems, interaction partners can change often, e.g. due to alterations in system structure or inclusion of new agents. Since the agents' benevolence cannot be assumed, they might not be willing to communicate their true intentions [27]. To deal with this situation, a reputation system can be used which combines the opinions of agents and generates recommendations [7]. This enables cooperation between agents that do not know or have only little experience with each other. To make adequate decisions, agents can rely on a combination of direct trust and reputation. To this end, several approaches [15, OCT5] propose to use confidence or similar metrics to dynamically weigh the influence of direct trust and reputation, e.g. depending on the number of direct experiences. Due to the subjective nature of trust and because agents might lie about the trustworthiness of others, it is often also desirable to weigh the impact a recommending agent, called *witness*, has on the reputation value. The *neighbour trust metric* [OCT6] as well as

*DTMAS* [28] propose to increase the influence of a witness with the similarity of the provided valuation to the one of the requesting agent. If the difference is too large, the witness can even be excluded from the calculation. This allows the system to deal with false reports. Further approaches that incentivise agents to provide truthful reports are discussed in Sect. 1.3.3. Providing reputation data can also be regarded as a special context in which witnesses are assessed according to the quality of their recommendations. In an even more fine-grained system, the context can also include information for which kind of interaction the recommendation is given. Whenever a reputation system is used, there has to be a consensus among the agents about the meaning of trust and reputation values. A common trust model can fulfil this purpose.

*Accountability, Deceit, and Collusion.* Open systems with little control over the agents are prone to exploitation from egoistic or malevolent agents. Therefore, special measures have to be taken to provide accountability of the agents and to prevent collusion. For an overview of attacks on trust and reputation management systems, see, e.g. [29]. Specific countermeasures are often system- or domain-specific, such as those presented for mobile ad-hoc networks in [30] or electronic markets in [31]. An important part of fraud prevention is a well-designed incentive system in combination with efficient monitoring facilities [32].

### 1.3 Different Uses of Trust in Open Self-Organising Systems

As discussed in Sect. 1.1, trust is traditionally used for selecting appropriate interaction partners. Bernard et al. [OCT7] call an agent's set of preferred interaction partners whose trust value is above a predefined threshold its *Implicit Trusted Community* (iTC). From the local view of a single agent, its interaction partners are selected through an implicit formation process. Note that this process is fully decentralised and thus not governed or controlled by an explicit authority. Because the agents do not coordinate their selections, the members of an iTC do not necessarily mutually trust each other. Yet this simple approach successfully excludes notoriously untrustworthy agents from most interactions.

In the following, we give an overview of four different uses of trust that extend this traditional use. First, we consider the trust-based formation of explicit organisations that allow large-scale open systems to deal with untrustworthy agents (see Sect. 1.3.1). Second, robust task and resource allocation promises to improve the system's stability and efficiency in uncertain environments (see Sect. 1.3.2). Third, uncertainties resulting from intentional misbehaviour can be reduced by means of appropriate incentives – employing trust as a sanctioning mechanism is one of several possibilities (see Sect. 1.3.3). Fourth, we outline measures how user trust in open environments can be increased (see Sect. 1.3.4).

### 1.3.1 *Trust to Structure Large-Scale Open Systems*

In essence, self-organisation enables a system to autonomously form and adapt a structure that supports its objectives under changing conditions. The main reasons for agents to form organisations are to achieve scalability and promote cooperation in order to accomplish their own or the system's goals [33]. While scalability is the result of the accompanying problem decomposition, cooperation is necessary due to the agents' limited resources and capabilities. There are a multitude of paradigms and algorithms for establishing organisations in literature, such as *teams* [33] and *coalition formation* [34]: While teams assume altruistic behaviour, coalition formation is used in systems consisting of self-interested and individually rational agents.

The participants of open systems might not only show self-interested behaviour but also lie about their capabilities, the utility of performing an action, etc. Consequently, the selection of suitable cooperation partners becomes even more important. Since suitable coalition structures depend on the agents' promised contributions, the system has to make sure that these promises are kept and all coalition members pursue a common goal. To this end, extensions of coalition formation incorporating trust into the agents' decisions have been presented in [35, 36]. In contrast to coalitions, *clans* [37] are long-lived. Given that cooperation is likely to be most beneficial and least uncertain with trustworthy agents, clans are groups of agents that mutually trust each other. A similar concept, called *Explicit Trusted Communities* (eTCs), for the domain of Desktop Grid Computing has been proposed in [OCT8]. The main difference to clans and coalitions is that each eTC is represented by an explicit manager which administrates memberships, deals with conflicts, and governs the participating agents with norms. By preferring interactions with trustworthy agents (or even restricting them to these agents), clans and eTCs incentivise untrustworthy agents to change their behaviour (see Sect. 1.3.3 for incentive mechanisms and norms). Ultimately, this procedure aims at a more efficient and robust system – at least with regard to the members of clans or eTCs. While these types of organisations are not necessarily limited to intentional misbehaviour, they assume that agents can be excluded from other parts of the system without jeopardising the overall system's stability and efficiency. This is why trustworthy agents can form exclusive groups.

However, there are situations in which untrustworthy agents can or should not be excluded from the system, e.g. if the system depends on their resources or if they can provide them in a particularly cost-efficient way. In power management systems, for instance, although the output of solar power plants is difficult to predict (their volatility is mirrored in low trust values), they should not be turned off because of their low-cost generation. If, in such a situation, scalability requires the agents to self-organise into subsystems, other types of organisations are needed to deal with untrustworthy agents. One possibility is the formation of *homogeneous partitionings* [OCT9] where organisations are as similar as possible with respect to certain criteria that have been identified as supporting the system's goals (including

their mean trustworthiness). This idea is based on the assumption that a centralised system imposes an upper bound on the ratio between trustworthy and untrustworthy agents: Given the uncertainties introduced by untrustworthy agents, the centralised control over trustworthy agents allows the system to fulfil its task as well as possible. If all organisations exhibit similar characteristics with respect to the identified criteria, such as a similar ratio between trustworthy and untrustworthy agents, they approximate the corresponding ratio of the centralised system. Consequently, they also inherit its positive properties. Ideally, this results in an organisational structure in which each organisation can deal with its untrustworthy agents internally without affecting or involving other organisations. In such situations, homogeneous partitioning increases the system's robustness and efficiency, and should be preferred to organisations consisting of homogeneous agents. A similar goal has been pursued in [38] where agents mitigate uncertainties originating from unintentional misbehaviour by forming coalitions in a way that they cancel each other out.

### ***1.3.2 Trust as a Basis for Robust Task or Resource Allocation***

In many applications, a MAS has to solve a task or resource allocation problem in which a set of tasks is to be allocated to agents, or a set of the agents have to provide a certain amount of resources in order to satisfy a given demand [39]. Due to the agents' limited resources and knowledge, they usually have to cooperate in order to achieve the goal. In open systems, finding an adequate allocation is even more difficult since agents might not provide resources or fulfil the task as promised and the actual demand that has to be satisfied or the resources required to perform a task might not be known exactly beforehand. Both types of uncertainties can be attributed to unintentional or intentional misbehaviour of the system's participants or its environment [OCT1]. If the system's stability or efficiency hinges on how well the agents fulfil the tasks or meet the demand – e.g. think about the demand of electric load in a smart grid application – techniques for robust task or resource allocation have to be regarded. In general, the way a robust allocation can be obtained depends on the type of misbehaviour.

*Unintentional misbehaviour* is introduced by external forces, such as current weather conditions. While this type of misbehaviour cannot be actively reduced, trust can be used to quantify and anticipate the uncertainties [10]. Incorporating trust into the decision-making process allows the system to optimise for *expectations*, such as the expected probability of success [40]. In [OCT10, OCT11], a self-organising middleware incorporating a trust-aware load-balancing mechanism assigns important services to trustworthy nodes in order to increase the services' expected availability. Similarly, participants of a Desktop Grid Computing system delegate the calculation of jobs to trustworthy agents, i.e. to members of their eTC, to improve their expected outcome (see Sect. 1.3.1). If the *predictability* (cf. “confidence”) of an agent's behaviour depends on its state, allocations can also be made in a way that promotes predictable behaviour [OCT1]. For highly



volatile environments in which dependencies in a sequence of observed behaviour have to be captured, a more expressive trust model, called *Trust-Based Scenario Trees* (TBSTs), has been proposed in [OCT2]. Basically, each TBST represents an empirical probability mass function that approximates the observed stochastic process. In contrast to trust models that capture the expected uncertainty or its variation, a TBST holds multiple possible scenarios, each with a probability of occurrence, of how the uncertainty might develop over a sequence of time steps. As opposed to the concept of *scenario trees* as known from the domain of operations research [41], TBSTs make only few assumptions about the underlying stochastic process. Further, they have been developed with the purpose of being *learned online* by agents with possibly low computational power. Combined with the principle of *stochastic programming* [42], agents can obtain robust allocations dynamically at runtime.

*Intentional misbehaviour* can be ascribed to agents that lie about some private information needed to decide about an adequate allocation, such as the cost or probability of performing a task successfully [40, 43]. Contrary to unintentional misbehaviour, uncertainties originating from intentional misbehaviour can be avoided. The field of *mechanism design* [40] studies how a system has to work in order to *incentivise* its self-interested, strategic, and individually rational participants to tell the truth. Further details concerning this matter are discussed in the following section.

### 1.3.3 *Trust as a Sanctioning and Incentive Mechanism*

Employing the techniques of *mechanism design* (MD) can guarantee efficiency (maximisation of the agents' overall utility), individual rationality (the agents' utility of participating in the scheme is non-negative), and incentive compatibility (the agents are best off revealing their true type) [44]. The latter property is of particular interest in open systems when agents have to be incentivised to disclose their private information needed to make decisions. In other words, MD can be used to incentivise individually rational agents to behave benevolently, that is, to ensure their trustworthy behaviour. *Fault-Tolerant MD* [43] and *Trust-Based MD* [40] address the issue of agents that have a probability of failure – quantified by a trust value – when performing an assigned task. Both approaches investigate the problem that reasonable task allocations depend on truthfully reported trust values. While each agent calculates and reports its own trust value in Fault-Tolerant MD [43], reputation values stemming from subjective trust measurements are considered in Trust-Based MD [40]. The ideas of MD have been adopted in various market-based approaches in which pricing mechanisms prevent agents from gaming the system [38, 45]. Depending on the regarded problem, it is often hard to devise a proper mechanism guaranteeing incentive compatibility, though, especially in case of unintentional misbehaviour. In these cases, it is still possible to use penalty schemes to increase the agents' risk that providing false reports or promises that

cannot be kept is detrimental to their utilities [44, OCT1]. Often, corresponding incentives can rely on the agents' trustworthiness. In electronic markets, trustworthy agents can obtain price premiums or price discounts [6]. In [OCT1], for instance, agents showing well-predictable behaviour can demand higher payments. Preferring trustworthy interaction partners or creating groups of trustworthy agents that benefit from a mutual increase in efficiency (cf. eTCs discussed in Sect. 1.3.1) also incentivises benevolent behaviour. These examples illustrate that trust in the sense of benevolent behaviour yields and, at the same time, embodies a form of *social capital* [46].

While the rules employed in these mechanisms are created at design time, open systems often have to be able to define, adjust, and implement behavioural guidelines in response to environmental and internal conditions at runtime. Such an adaptability is akin to Ostrom's principle of "congruence" that states that sustainable management of commons requires to "match rules governing use of common goods to local needs and conditions" [47]. While stemming from economic and sociological research, these Ostrom's principles have been recognised as the foundations for self-organising electronic institutions as well [48]. In *normative MAS* [49], *normative institutions* enact and enforce *norms* [50] to influence the agents' behaviour indirectly. Each norm describes a behavioural rule and a sanction that is imposed if the rule is not followed. A sanction might be punitive fines or a (temporary) reduction of the violator's reputation value. The latter type of sanction treats reputation in the sense of social capital such that its reduction incentivises trustworthy behaviour in the long run. If an agent did not violate a norm on purpose, if it compensates for the violation, or if the violation was inevitable, the institution might also abstain from a sanction, which introduces a form of *forgiveness* [51, 52]. Essentially, norms have to contribute to reaching the system's goal. In eTCs (see Sect. 1.3.1), managers take on the role of normative institutions. If a manager detects an attack, it defends its community by adjusting the set of norms, e.g. by regulating the delegation and the acceptance of jobs in case of a trust breakdown – a situation in which even the reputation of benevolent agents declines [OCT12]. To enforce norms, an institution must not only be able to react with sanctions but also to detect their violation. Since monitoring an agent's behaviour comes at a price, Edenhofer et al. [OCT13] proposed to couple the effort put into surveillance to the number of received accusations. Especially when regarding trust as the basis of delegation [18], norms can also be understood as social laws governing the delegation of institutional power [53]. In this case, norms represent explicit permissions that have to be acquired before a specific action may be performed.

### ***1.3.4 Increasing User Trust in Open Environments***

Beyond the use of trust to qualify the relationships between software agents (cf. system-to-system trust in Sect. 1.1), it can also be applied to describe the social relationships between the users and the system (cf. user-to-system trust). Recent

advances in sensor technologies and context recognition enable us to capture the users' physical context continuously and to personalise information and services to them in real-time. Apart from simply providing information, context-aware systems can also allow users to manipulate or share data or even act autonomously on their behalf. Combined with advances in display and wireless technologies, users can employ these systems basically anytime and anywhere. While these so-called ubiquitous environments offer great benefits to users, they also raise a number of challenges. In particular, they might show a behaviour that negatively affects user trust. Examples include (1) highly dynamic situations where the rationale behind the system's actions is no longer apparent to the user [54], (2) implicit interactions through proxemic behaviour where the user no longer feels in control [55], or (3) privacy issues [56]. Hence, there is an enormous need for sophisticated trust management in ubiquitous environments in order to ensure that such environments will find acceptance among users.

While most work in the area of computational trust models aims to develop trust metrics that determine, on the basis of objective criteria, whether a system should be trusted or not, not much interest has been shown towards trust experienced by a user when interacting with a system. A system may be robust and secure, but nevertheless be perceived as not very trustworthy by a user, e.g. because its behaviour appears opaque or hard to control. Following the terminology by Castelfranchi and Falcone [57], a focus is put on the affective forms of trust that are based on the user's appraisal mechanisms. Therefore, the objective must be to develop a computational trust model that captures how a system – and more specifically a ubiquitous environment – is perceived by a user while interacting with it.

Many approaches found in literature aim to identify trust dimensions that influence the user's feeling of trust. This is an extension to the trust models as discussed in Sect. 1.2, even though facets of trust play a role in open self-organising systems as well [OCT14]. Trust dimensions that have been researched in the context of internet applications and e-commerce include reliability, dependability, honesty, truthfulness, security, competence, and timeliness, see, e.g. [58, 59]. Tschannen et al. [60], who are more interested in the sociological aspects of trust, introduce willing vulnerability, benevolence, reliability, competence, honesty, and openness as the constituting facets of trust, although their work does not focus on trust in software. Researchers working on adaptive user interfaces consider transparency as a major component of trust, see, e.g. [61]. Trust dimensions have formed the underlying basis of many conceptual models of trust. However, incorporating them into a computational model of trust is not a trivial task.

With the *User Trust Model* (UTM) [62], such a computational model of trust was introduced, along with a decision-theoretic approach to trust management for ubiquitous and self-adaptive environments. The UTM is based on Bayesian networks and, following ideas put forward by Yan et al. [63], assesses the users' trust in a system, monitors it over time, and applies appropriate system reactions to maintain users' trust in critical situations. In a smart office application, for example, the system could automatically switch off the lights because it senses that it is

bright enough outside, but might actually decide against it if it assesses that such an action would have a negative impact on the user's trust due to a lack of control and transparency.

## 1.4 Conclusion

The potential of computational trust in open self-organising systems is substantial. As we outlined in this chapter, trust models can increase a system's fitness by providing a means to optimise for the most likely or most risky future states; they can decrease information asymmetry; they can be used in combination with sanctioning and incentive mechanisms in normative frameworks codifying behavioural guidelines; and they can enable the formation of a system structure supporting the functions of the system optimally. If trust models are used to represent the social relationships of a system's users [OCT15], the system can, for instance, even make robust decisions with regard to the users' privacy. The basic principles are the same for all of the uses shown here. They can all be applied on the basis of an understanding of trust that puts the concepts of interactions, contracts, and experiences at its core and is compatible with many trust models available in the literature.

In all cases, trust increases the efficiency and robustness of open self-organising systems by mitigating uncertainties originating from a system's unknown participants and the environment it is exposed to.

**Acknowledgements** This research is partly sponsored by the research unit *OC-Trust* (FOR 1085) of the German Research Foundation.

## References

1. Pasquier, P., Flores, R., Chaib-draa, B.: Modelling flexible social commitments and their enforcement. In: Gleizes, M.-P., Omicini, A., Zambonelli, F. (eds.) *Engineering Societies in the Agents World V*, vol. 3451, pp. 898–898. Springer, Berlin/Heidelberg (2005). ISBN:978-3-540-27330-1
2. Artikis, A., Pitt, J.: Specifying open agent systems: a survey. In: Artikis, A., Picard, G., Vercouter, L. (eds.) *Engineering Societies in the Agents World IX*, vol. 5485, pp. 29–45. Springer, Berlin/Heidelberg (2009). ISBN:978-3-642-02561-7
3. Cossentino, M., Gaud, N., Hilaire, V., Galland, S., Koukam, A.: ASPECS: an agent-oriented software process for engineering complex systems. *Auton. Agents Multi-agent Syst.* **20**, 260–304 (2010).
4. Boon, S., Holmes, J.: The dynamics of interpersonal trust: resolving uncertainty in the face of risk. In: Hinde, R., Groebel, J. (eds.) *Cooperation and Prosocial Behaviour*, pp. 190–211. Cambridge University Press, Cambridge (1991)
5. Rousseau, D., Sitkin, S., Burt, R., Camerer, C.: Not so different after all: a cross-discipline view of trust. *Acad. Manag. Rev.* **23**, 393–404 (1998)

6. Ba, S., Pavlou, P.: Evidence of the effect of trust building technology in electronic markets: price premiums and buyer behavior. *MIS Q.* **26**, 243–268 (2002)
7. Mui, L., Mohtashemi, M., Halberstadt, A.: A computational model of trust and reputation. In: *Proceedings of the 35th Hawaii International Conference on System Sciences (HICSS'02)*, Big Island, pp. 188–196 (2002)
8. Corritore, C., Kracher, B., Wiedenbeck, S.: On-line trust: concepts, evolving themes, a model. *Int. J. Hum.-Comput. Stud.* **58**, 737–758 (2003)
9. Koller, M.: Risk as a determinant of trust. *Basic Appl. Soc. Psychol.* **9**, 265–276 (1988)
10. Ramchurn, S., Huynh, D., Jennings, N.: Trust in multi-agent systems. *Knowl. Eng. Rev.* **19**, 1–25 (2004)
11. Marsh, S.P.: Formalising trust as a computational concept. PhD thesis, University of Stirling Digital Repository (1994)
12. Kantert, J., Edenhofer, S., Tomforde, S., Hähner, J., Müller-Schloer, C.: Normative control – controlling open distributed systems with autonomous entities. In: Reif, W., Anders, G., Seebach, H., Steghöfer, J.-P., André, E., Hähner, J., Müller-Schloer, C., Ungerer, T. (eds.) *Autonomic Systems*, vol. 7, pp. 87–123 (2016)
13. Mayer, R.C., Davis, J.H., Schoorman, F.D.: An integrative model of organizational trust (English). *Acad. Manag. Rev.* **20**, 709–734 (1995). ISSN:03637425
14. Yu, H., Shen, Z., Leung, C., Miao, C., Lesser, V.: A survey of multi-agent trust management systems. *IEEE Access* **1**, 35–50 (2013)
15. Sabater, J., Sierra, C.: Social regret, a reputation model based on social relations. *ACM SIGecom Exch.* **3**, 44–56 (2001)
16. Ramchurn, S.D., Jennings, N.R., Sierra, C., Godo, L.: Devising a trust model for multiagent interactions using confidence and reputation. *Appl. Artif. Intell.* **18**, 833–852 (2004)
17. Jonker, C., Treur, J.: Formal analysis of models for the dynamics of trust based on experiences. In: Garijo, F., Boman, M. (eds.) *Multi-agent System Engineering*, vol. 1647, pp. 221–231. Springer, Berlin/Heidelberg (1999). ISBN:978-3-540-66281-5
18. Falcone, R., Castelfranchi, C.: Social trust: a cognitive approach. In: Castelfranchi, C., Tan, Y.-H. (eds.) *Trust and Deception in Virtual Societies*, pp. 55–90. Kluwer Academic, Norwell (2001). ISBN:0-7923-6919-X
19. McKnight, D., Cummings, L., Chervany, N.: Initial trust formation in new organizational relationships. *Acad. Manag. Rev.* **23**, 473–490 (1998)
20. Jøsang, A., Pope, S.: Semantic constraints for trust transitivity. In: *Proceedings of the 2nd Asia-Pacific Conference on Conceptual Modelling*, vol. 43, pp. 59–68. Australian Computer Society, Newcastle (2005). ISBN:1-920-68225-2
21. Good, D.: Individuals, interpersonal relations, and trust. In: Gambetta, D. (ed.) *Trust: Making and Breaking Cooperative Relations*, pp. 31–48. Department of Sociology, University of Oxford (2000)
22. Jones, K.: Trust as an affective attitude (English). *Ethics* **107**, 4–25 (1996). ISSN:00141704
23. He, R., Niu, J., Zhang, G.: CBTM: A trust model with uncertainty quantification and reasoning for pervasive computing. In: Pan, Y., Chen, D., Guo, M., Cao, J., Dongarra, J. (eds.) *Parallel and Distributed Processing and Applications*, pp. 541–552. Springer, Berlin/Heidelberg (2005). ISBN:978-3-540-29769-7
24. Wang, Y., Singh, M.P.: Formal trust model for multiagent systems. In: *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, pp. 1551–1556. Morgan Kaufmann Publishers Inc., San Francisco (2007)
25. Kuter, U., Golbeck, J.: Using probabilistic confidence models for trust inference in web-based social networks. *ACM Trans. Internet Technol.* **10**, 8:1–8:23 (2010). ISSN:1533-5399
26. Marsh, S.: Optimism and pessimism in trust. In: *Proceedings of the Ibero-American Conference on Artificial Intelligence (IBERAMIA '94)*, Caracas. McGraw-Hill Publishing (1994)
27. Schillo, M., Funk, P., Rovatsos, M.: Using trust for detecting deceitful agents in artificial societies. *Appl. Artif. Intell.* **14**, 825–848 (2000)

28. Aref, A.M., Tran, T.T.: A decentralized trustworthiness estimation model for open, multiagent systems (DTMAS). *J. Trust Manag.* **2**, 1–20 (2015)
29. Jøsang, A.: Robustness of trust and reputation systems: does it matter? In: Proceedings of IFIPTM International Conference on Trust Management (IFIPTM 2012), Surat. Springer (2012)
30. Sun, Y., Han, Z., Liu, K.: Defense of trust management vulnerabilities in distributed networks. *IEEE Commun. Mag.* **46**, 112–119 (2008). ISSN:0163-6804
31. Yao, Y., Ruohomaa, S., Xu, F.: Addressing common vulnerabilities of reputation systems for electronic commerce. *J. Theor. Appl. Electron. Commer. Res.* **7**, 1–20 (2012)
32. Grossi, D., Aldewereld, H., Dignum, F.: Ubi Lex, Ibi Poena: designing norm enforcement in E-institutions. In: Noriega, P., Vázquez-Salceda, J., Boella, G., Boissier, O., Dignum, V., Fornara, N., Matson, E. (eds.) *Coordination, Organizations, Institutions, and Norms in Agent Systems II*, vol. 4386, pp. 101–114. Springer, Berlin/Heidelberg (2007). ISBN:978-3-540-74457-3
33. Horling, B., Lesser, V.: A survey of multi-agent organizational paradigms. *Knowl. Eng. Rev.* **19**, 281–316 (2004)
34. Shehory, O., Kraus, S.: Methods for task allocation via agent coalition formation. *Artif. Intell.* **101**, 165–200 (1998)
35. Breban, S., Vassileva, J.: Long-term coalitions for the electronic marketplace. In: Proceedings of the E-Commerce Applications Workshop at the Canadian AI Conference, Ottawa (2001)
36. Breban, S., Vassileva, J.: A coalition formation mechanism based on inter-agent trust relationships. In: Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems: Part 1, pp. 306–307. ACM, New York (2002). ISBN:1-58113-480-0
37. Griffiths, N., Luck, M.: Coalition formation through motivation and trust. In: Proceedings of the 2nd International Joint Conference on Autonomous Agents and Multiagent Systems, pp. 17–24. ACM, Melbourne (2003). ISBN:1-58113-683-8
38. Chalkiadakis, G., Robu, V., Kota, R., Rogers, A., Jennings, N.R.: Cooperatives of distributed energy resources for efficient virtual power plants. In: Proceedings of the 10th International Conference on Autonomous Agents and Multiagent Systems. International Foundation for Autonomous Agents and Multiagent Systems, Taipei, vol. 2, pp. 787–794 (2011). ISBN:0-9826571-6-1, 978-0-9826571-6-4
39. Chevaleyre, Y., Dunne, P.E., Endriss, U., Lang, J., Lemaître, M., Maudet, N., Padget, J., Phelps, S., Rodríguez-aguilar, J.A., Sousa, P.: Issues in multiagent resource allocation. *Informatica* **30**, 3–31 (2006)
40. Dash, R.K., Ramchurn, S.D., Jennings, N.R.: Trust-based mechanism design. In: Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems, vol. 2, pp. 748–755. IEEE Computer Society, Washington, DC (2004)
41. Shapiro, A., Dentcheva, D., Ruszczyński, A.: *Lectures on Stochastic Programming: Modeling and Theory*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia (2014)
42. Hentenryck, P.V., Bent, R.: *Online Stochastic Combinatorial Optimization*. MIT Press, Cambridge/London (2009)
43. Porter, R., Ronen, A., Shoham, Y., Tennenholtz, M.: Mechanism design with execution uncertainty. In: Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence, Edmonton, pp. 414–421 (2002)
44. Dash, R., Vytelingum, P., Rogers, A., David, E., Jennings, N.: Market-based task allocation mechanisms for limited-capacity suppliers. *Syst. Man Cybern. Part A: IEEE Trans. Syst. Hum.* **37**, 391–405 (2007)
45. Vytelingum, P., Voice, T., Ramchurn, S., Rogers, A., Jennings, N.: Agent-based micro-storage management for the smart grid. In: Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems, Toronto, vol. 1, pp. 39–46 (2010)
46. Pitt, J., Nowak, A.: The reinvention of social capital for socio-technical systems [Special Section Introduction]. *IEEE Technol. Soc. Mag.* **33**, 27–80 (2014)
47. Ostrom, E.: *Governing the Commons: The Evolution of Institutions for Collective Action*. Cambridge University Press, Cambridge/New York (1990)

48. Pitt, J., Schaumeier, J., Artikis, A.: Axiomatization of socio-economic principles for self-organizing institutions: concepts, experiments and challenges. *ACM Trans. Auton. Adapt. Syst. (TAAS)* **7**, 39 (2012)
49. Boella, G., Pigozzi, G., van der Torre, L.: Normative systems in computer science – ten guidelines for normative multiagent systems. In: Boella, G., Noriega, P., Pigozzi, G., Verhagen, H. (eds.) *Normative Multi-agent Systems*. Schloss Dagstuhl, Leibniz-Zentrum fuer Informatik (2009)
50. Conte, R., Castelfranchi, C.: Norms as mental objects. From normative beliefs to normative goals. In: Castelfranchi, C., Müller, J.-P. (eds.) *From Reaction to Cognition*, vol. 957, pp. 186–196. Springer, Berlin/Heidelberg (1995). ISBN:978-3-540-60155-5
51. Vasalou, A., Pitt, J.: Reinventing forgiveness: a formal investigation of moral facilitation. In: Herrmann, P., Issarny, V., Shiu, S. (eds.) *Trust Management*, vol. 3477, pp. 39–90. Springer, Berlin/Heidelberg (2005). ISBN:978-3-540-26042-4
52. Marsh, S., Briggs, P.: Examining trust, forgiveness and regret as computational concepts. In: Golbeck, J. (ed.) *Computing with Social Trust*, pp. 9–43. Springer, London (2009). ISBN:978-1-84800-356-9
53. Artikis, A., Sergot, M., Pitt, J.: Specifying norm-governed computational societies. *ACM Trans. Comput. Log. (TOCL)* **10**, 1 (2009)
54. Rothrock, L., Koubek, R., Fuchs, F., Haas, M., Salvendy, G.: Review and reappraisal of adaptive interfaces: toward biologically inspired paradigms. *Theor. Issues Ergon. Sci.* **3**, 47–84 (2002)
55. Müller, J., Exeler, J., Buzeck, M., Krüger, A.: ReflectiveSigns: digital signs that adapt to audience attention. In: *Proceedings of 7th International Conference on Pervasive Computing*, pp. 17–24. Springer, Berlin/Heidelberg (2009)
56. Röcker, C., Hinske, S., Magerkurth, C.: Intelligent privacy support for large public displays. In: *Proceedings of Human-Computer Interaction International 2007 (HCI'07)*. Beijing, China (2007)
57. Castelfranchi, C., Falcone, R.: *Trust Theory: A Socio-Cognitive and Computational Model*. Wiley, Hoboken (2010)
58. Grandison, T., Sloman, M.: A survey of trust in internet applications. *IEEE Commun. Surv. Tutor.* **3**, 2–16 (2000)
59. Kini, A., Choobineh, J.: Trust in electronic commerce: definition and theoretical considerations. *Proc. Hawaii Int. Conf. Syst. Sci.* **31**, 51–61 (1998)
60. Tschannen-Moran, M., Hoy, W.: A multidisciplinary analysis of the nature, meaning, and measurement of trust. *Rev. Educ. Res.* **70**, 547 (2000)
61. Glass, A., McGuinness, D.L., Wolverton, M.: Toward establishing trust in adaptive agents. In: *Proceedings of the 13th International Conference on Intelligent User Interfaces (IUI '08)*, pp. 227–236. ACM, New York (2008)
62. Hammer, S., Wißner, M., André, E.: A user trust model for automatic decision-making in ubiquitous and self-adaptive environments. In: Reif, W., Anders, G., Seebach, H., Steghöfer, J.-P., André, E., Hähner, J., Müller-Schloer, C., Ungerer, T. (eds.) *Autonomic Systems*, vol. 7, pp. 55–86 (2016)
63. Yan, Z., Holtmanns, S.: *Computer Security Privacy and Politics: Current Issues, Challenges and Solutions*, pp. 290–323. IGI Global, Hershey, USA (2008)

## References Originating from the OC-Trust Project

- OCT1. Anders, G., Schiendorfer, A., Siefert, F., Steghöfer, J.-P., Reif, W.: Cooperative resource allocation in open systems of systems. *ACM Trans. Auton. Adapt. Syst.* **10**, 11:1–11:44 (2015)

- OCT2. Anders, G., Siefert, F., Steghöfer, J.-P., Reif, W.: Trust-based scenarios – predicting future agent behavior in open self-organizing systems. In: Elmenreich, W., Dressler, F., Loreto, V. (eds.) *Self-Organizing Systems*, vol. 8221, pp. 90–102. Springer, Berlin/Heidelberg (2014). ISBN:978-3-642-54139-1
- OCT3. Steghöfer, J.-P., Reif, W.: Die Guten, die Bösen und die Vertrauenswürdigen–Vertrauen im Organic Computing. *Informatik-Spektrum* **35**, 119–131 (2012) (in German)
- OCT4. Kiefhaber, R., Anders, G., Siefert, F., Ungerer, T., Reif, W.: Confidence as a means to assess the accuracy of trust values. In: *Proceedings of the 11th IEEE International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom-2012)*, Liverpool, pp. 690–697. IEEE (2012)
- OCT5. Kiefhaber, R., Jahr, R., Msadek, N., Ungerer, T.: Ranking of direct trust, confidence, and reputation in an abstract system with unreliable components. In: *Ubiquitous Intelligence and Computing, 2013 IEEE 10th International Conference on and 10th International Conference on Autonomic and Trusted Computing (UIC/ATC)*, Sorrento Peninsula, Italy, pp. 388–395 (2013)
- OCT6. Kiefhaber, R., Hammer, S., Savs, B., Schmitt, J., Roth, M., Kluge, F., Andre, E., Ungerer, T.: The neighbor-trust metric to measure reputation in organic computing systems. In: *2011 Fifth IEEE Conference on Self-Adaptive and Self-organizing Systems Workshops (SASOW)*, Ann Arbor, pp. 41–46 (2011)
- OCT7. Bernard, Y., Klejnowski, L., Hähner, J., Müller-Schloer, C.: Towards trust in desktop grid systems. In: *IEEE International Symposium on Cluster Computing and the Grid*, pp. 637–642. IEEE Computer Society, Los Alamitos (2010). ISBN:978-0-7695-4039-9
- OCT8. Klejnowski, L.: *Trusted community: a novel multiagent organisation for open distributed systems*. PhD thesis, Leibniz Universität Hannover (2014). <http://edok01.tib.uni-hannover.de/edoks/e01dh11/668667427.pdf>
- OCT9. Anders, G., Siefert, F., Reif, W.: A Heuristic for Constrained Set Partitioning in the Light of Heterogeneous Objectives. In: *Agents and Artificial Intelligence*. LNAI. Lisbon, Portugal (2015)
- OCT10. Msadek, N., Kiefhaber, R., Ungerer, T.: A trustworthy fault-tolerant and scalable self-configuration algorithm for organic computing systems. *J. Syst. Archit.* **61**(10), 511–519 (2015)
- OCT11. Msadek, N., Kiefhaber, R., Ungerer, T.: Trustworthy self-optimization in organic computing environments. *Architecture of Computing Systems – ARCS 2015*, pp. 123–134. Porto, Portugal (2015)
- OCT12. Kantert, J., Scharf, H., Edenhofer, S., Tomforde, S., Hähner, J., Müller-Schloer, C.: A graph analysis approach to detect attacks in multi-agent-systems at runtime. In: *2014 IEEE Eighth International Conference on Self-Adaptive and Self-Organizing Systems*, pp. 80–89. IEEE, London (2014)
- OCT13. Edenhofer, S., Stifter, C., Jänen, U., Kantert, J., Tomforde, S., Hähner, J., Müller-Schloer, C.: An accusation-based strategy to handle undesirable behaviour in multi-agent systems. In: *2015 IEEE Eighth International Conference on Autonomic Computing Workshops (ICACW)*. Grenoble, France (2015)
- OCT14. Steghöfer, J.-P., Kiefhaber, R., Leichtenstern, K., Bernard, Y., Klejnowski, L., Reif, W., Ungerer, T., André, E., Hähner, J., Müller-Schloer, C.: Trustworthy organic computing systems: challenges and perspectives. In: Xie, B., Branke, J., Sadjadi, S., Zhang, D., Zhou, X. (eds.) *Autonomic and Trusted Computing*, vol. 6407, pp. 62–76. Springer, Berlin/Heidelberg (2010). ISBN:978-3-642-16575-7
- OCT15. Kurdyukova, E., Bee, K., André, E.: Friend or foe? Relationship-based adaptation on public displays. In: *Proceedings of the Second International Conference on Ambient Intelligence*, pp. 228–237. Springer, Amsterdam (2011). ISBN:978-3-642-25166-5