

1 **Chapter 31**
2 **Commitments in multiagent systems**
3 **Some history, some confusions,**
4 **some controversies, some prospects**

5 Munindar P. Singh

6 **Abstract** The notion of commitments as a foundation for understanding in-
7 teractions among agents has been under development for about twenty years.
8 Cristiano Castelfranchi has contributed to clarifying the conception of com-
9 mitments by bringing in insights from social psychology. In this essay, I briefly
10 review the conceptual development of commitments in multiagent systems,
11 identifying the key themes and some lingering confusions. I also highlight
12 some ongoing debates with Castelfranchi and some promising directions for
13 future research.

14 **1 Introduction**

15 Cristiano Castelfranchi writes about agents like Michelangelo painted his fres-
16 coes. No, I don't mean to suggest that Cristiano writes lying precariously on
17 his back on scaffolding twenty meters above the floor—although one can never
18 be too sure about the ways of Italian intellectuals. Seriously, though, I do mean
19 to suggest that Cristiano naturally envisions and describes complex scenes
20 with many characters and details. The effect is beautiful indeed.

21 My goal in this short essay, by contrast, is to come to these frescoes as a
22 computer scientist—generally, focusing on a few characters and their particular
23 details in an attempt to understand some components of the scene better.

24 Professor Castelfranchi has made varied and numerous contributions to
25 identifying, developing, and popularizing the social perspective on multi-
26 agent systems. Specifically, I want to focus on the notion of commitments,
27 which Professor Castelfranchi and I have been contemplating and discussing
28 for nearly two decades [3, 4, 38]. This is not to suggest that others haven't
29 contributed to this topic—the study of commitments has become a veritable

Munindar P. Singh
Department of Computer Science, North Carolina State University, USA
e-mail: singh@ncsu.edu

1 cottage industry—but merely to focus the presentation on themes that interest
2 Professor Castelfranchi and me the most.

3 2 A Brief Retrospective

4 Commitments in multiagent systems turn out to be quite different from “com-
5 mitments” as have long been discussed in artificial intelligence (AI) and phi-
6 losophy. In traditional AI work, a commitment was understood as the com-
7 mitment of a single agent to some belief or to some course of action. For
8 example, the AI planning literature of the 1970s advocated an approach called
9 least-commitment planning [35] wherein a planner (working as part of or on
10 behalf of a single agent) would create a plan that left as many of the options
11 of the agent open as it could and for as long as it could. That’s a fine idea for a
12 single-agent setting. Notice—as an aside—that in such a setting a commitment is
13 not quite desirable—an agent is best off when its commitments are minimized.

14 In the mid to late 1980s, when research began in earnest on multiagent sys-
15 tems, researchers adopted the notion of commitment as a way to understand
16 organizations of agents. A commitment in a multiagent system captures a re-
17 lationship between two parties. A traditional planning-style commitment to
18 one’s plans would not suffice. Even though the researchers recognized this,
19 partly because they came from an AI background, they came to the notion of
20 commitments with an attendant mentalist bias [30]. Thus they distinguished
21 multiagent commitments from planning commitments, but only to the extent
22 of somehow reducing multiagent commitments to combinations of mutual
23 beliefs and intentions. A mutual belief between two or more agents is a propo-
24 sition p where each believes p and each believes that each believes p , and so
25 on, to arbitrary nesting [8]. Mutual or joint intentions are similar in spirit,
26 though somewhat more subtle [30]. In other words, traditional researchers re-
27 tained their mentalist perspective but hoped that the mutuality of the beliefs
28 or intentions would provide the glue between the agents.

29 But, as Professor Castelfranchi has eloquently and forcefully argued, so-
30 cial relationships are irreducible to the mental attitudes. And especially in
31 multiagent systems we are concerned with the modeling and enactment of in-
32 teractions of autonomous and heterogeneous agents. Thus commitments can
33 easily exist or fail to exist with or without any beliefs or intentions on part of
34 any of the agents. I return to this topic in Section 4 along with addressing some
35 other confusions.

36 In contrast, the social and organizational metaphors provided a more
37 straightforward way to think of multiagent systems, and especially as a way
38 to formulate commitments. It has been long known that human organizations
39 develop and apply standard operating procedures—as, for example, explained
40 by Herbert Simon [37]. And, especially in settings where there may be no
41 mathematical guarantee of obtaining a rigorously correct state or outcome,
42 applying a standard operating procedure would be the rational way for an
43 organization to proceed—in essence, we would define the state or outcome

1 emerging from such a procedure as being correct. A member of the organi-
2 zation, when faced with a particular situation, could act according to any of
3 the applicable standard operating procedures. Even if the particular outcome
4 in that situation turned out to be undesirable or even harmful, the member
5 would generally not be considered as having been in violation. For example,
6 if a patient collapses in an apparent heart failure, a paramedic nurse may be
7 expected to give the patient a shot of nitrates. The nurse would be deemed
8 to have done the right thing if he gives the patient the specified amount of
9 the recommended medication even if he is unable to save the patient's life or
10 saves the patient's life but inadvertently causes other complications. Clearly
11 there are cases where the standard expectations may be higher and a member
12 of an organization would need to select an appropriate operating procedure in
13 order to avoid all blame. Further, the expectations can vary depending upon
14 the role and qualifications of the member involved. In the above example, we
15 may expect an emergency physician or a cardiologist to consider additional
16 information and potential risks beyond what we expect a nurse to consider in
17 deciding a course of action. But regardless of whether we consider a simplified
18 notion of an operating procedure or a more complex one, the common feature
19 is that the organization empowers its members to act in circumstances that are
20 far from ideal.

21 To me, the foregoing line of thought led to an inkling of an idea that an
22 organization be able to commit to a course of action. More pertinently, the
23 commitment here arose from the member to the organization. Such thinking
24 led me to distinguish two kinds of commitments: (1) an internal one, which I
25 then termed psychological or P-commitment and (2) an external one, which I
26 termed social or S-commitment [38]. Psychological commitment is the standard
27 concept in AI. Social commitment is the concept that we now refer to as
28 commitment in the multiagent systems community. The AI researchers resisted
29 social commitments. I am grateful to Professor Castelfranchi for lending his
30 support to this area when it was emerging.

31 Social commitments have some interesting features distinguishing them
32 from psychological commitments. First, a social commitment is directed from
33 one party (its *debtor*) to another (its *creditor*). This terminology reflects the
34 intuition that the debtor is committed to doing something for the creditor.

35 The idea was to distinguish this from the more obvious notion of a benefi-
36 ciary. Specifically, a commitment may be directed toward one party but the
37 beneficiary might be another. For example, a shipper may commit to a mer-
38 chant to deliver a package to a customer. Here the shipper would be the debtor
39 and the merchant the creditor. The apparent beneficiary, the recipient of the
40 package, may show up only within the body of the condition that the shipper
41 commits to bring about. Notice that the logical form of the above commitment
42 is the same as of the commitment where the local police constable commits
43 to the district attorney to deliver a subpoena to or arrest a citizen. We would
44 generally not think of the citizen being subject to a subpoena as being a benefi-
45 ciary. For this reason, it is advisable to leave the value judgments of who is
46 the beneficiary and who is not outside of the general concept of commitments.
47 Indeed, such value judgments are often accompanied by presumptions about

1 various psychological concepts, which too we ought to minimize in the general
2 theory.

3 The second interesting, and even less common, aspect of social commit-
4 ments was the idea of incorporating an *organizational context* into the notion of
5 commitments. The organizational context of a commitment describes the orga-
6 nization or “system” in which the commitment arises, providing support of the
7 normative backdrop for commitments and interactions among autonomous
8 parties. The debtor and creditor of the commitment would thus generally be
9 members of the context organization. An example would be a commitment
10 from a seller to a buyer operating within the eBay marketplace wherein the
11 seller is committed to shipping some goods. That commitment references eBay
12 as its organizational context. Here, eBay might penalize a seller who doesn’t
13 discharge the commitment.

14 A related intuition is that agents can be composed. In other words, what,
15 from one perspective, appears to be an individual entity and functions as
16 one entity (and therefore is a well-defined entity) might well, from a different
17 perspective, turn out to be internally structured. For instance, a corporation or a
18 university might function and interact as if it were an individual, for example,
19 by entering into contracts with others. Yet, from an internal perspective, it
20 would generally consist of several agents.

21 Combining the above intuition with the context-based view of commit-
22 ments is that it enables us to express complex domain structures in a simple
23 manner. For instance, we can imagine a team as the organizational context of
24 the several commitments that tie together its members. However, the team is
25 itself constructed from its members. Thus we would naturally model commit-
26 ments between the team viewed as an agent and each of its members. Such
27 commitments might capture the principal intuitions of teamwork such as that
28 a member of a team should support the other members of the team in succeed-
29 ing with their goals, and that the team-members should coordinate with one
30 another to accomplish their common goals. The specification of commitments
31 between the members and the team help codify such relationships. The team-
32 members need not form mutual beliefs or joint intentions with one another, as
33 traditional approaches require [30, 23], because the essence of the relationship
34 between them can be captured through the commitments. In particular, a team-
35 member may not even know who the other team-members are in order to form
36 a social relationship with them through the common identity of the team to
37 which they belong. A further benefit is that the relationships naturally express
38 the rules of encounter of the team and thus support the expectations that the
39 team-members might form on each other. Additionally, the relationship can
40 potentially be realized in a variety of ways. For example, the members of a
41 team may join it one by one and a team-member may leave and another may
42 join without altering the essential fabric of the team. Or, the team-members
43 may all join at once.

44 An important aspect of commitments is that they can be manipulated [40].
45 A debtor may create or cancel a commitment; a creditor may release it. More
46 interestingly, a debtor may delegate a commitment to a new debtor, and a
47 creditor may assign it to a new creditor. Such manipulations provide a high-

1 level and systematic way in which the social state can progress. Fornara and
2 Colombetti [19], [20] have studied the operationalization of commitments to
3 support such manipulations. With Xing [52] and Chopra and Desai [46], I have
4 further developed patterns involving the manipulation of commitments that
5 support useful properties.

6 This wasn't always emphasized in the early works, but the conditionality
7 of commitments is important. By default, commitments are conditional,
8 involving an antecedent and a consequent, and unconditional commitments
9 are merely the case where the antecedent is true. In logical terms, the conditionality
10 of commitments resembles that of a strong conditional rather than a
11 material conditional [43]. A commitment becomes and stays *detached* or *dis-*
12 *charged* when, respectively, its antecedent and consequent become true. There
13 is no presumption of temporal order between the detach and discharge of a
14 commitment. A commitment that is detached but fails to discharge indicates
15 a violation.

16 It is worth distinguishing two major kinds of commitments. *Practical com-*
17 *mitments*—as commonly seen in formalizations of business processes—are about
18 what the debtor would bring about. *Dialectical commitments*—as commonly seen
19 in formalizations of argumentation—are about what the debtor stakes a claim
20 on. The import of the two kinds of commitments is quite different and parallels
21 the distinction between goals and beliefs, respectively. Practical commitments
22 call for action and thus relate to present or future actions. Dialectical com-
23 mitments call for a condition holding and thus can relate to past, present, or
24 future.

25 That the two kinds of commitment are distinct has been known for years
26 and, in particular, finds discussion in some of Professor Castelfranchi's work
27 wherein he provides the clearest exposition of it. However, the distinction
28 seems to have been lost in the agents literature until recently. I have sought
29 to revive this distinction in conjunction with a proposed formal semantics for
30 commitments [43].

31 **3 What are Commitments Good For?**

32 In a nutshell, commitments form a key element, arguably the most important
33 element, of the social state of two or more interacting agents.

34 Commitments are important because they help us address the tradeoffs
35 between and reconcile the tension between autonomy and interdependence.
36 On the one hand, we would like to model our agents as being autonomous with
37 respect to each other. On the other hand, it is clear that if the agents were fully
38 autonomous, then we would have not a multiagent system in the true sense
39 of the term, but merely a number of agents that happen to coexist in a shared
40 environment. Such a system would exhibit no useful structure. Further, it is
41 clear that autonomous agents must be able to cooperate and compete with each
42 other, and carry out complex interactions. If there were no interdependence,
43 the agents would be nearly useless. Professor Castelfranchi and his colleagues

1 first articulated the importance of such interdependence among agents and
2 explored varieties of it [36]. Similar intuitions and elaborating the connection
3 with autonomy arise in newer work [27]. Commitments provide a natural way
4 to characterize the bounds of autonomy and interdependence without getting
5 bogged down in low-level details.

6 *3.1 Commitments for Business Protocols*

7 A business protocol characterizes how a family of interactions involving two
8 or more business partners may proceed. What makes a business protocol
9 “business” is that the interactions it characterizes involve business relation-
10 ships. The classic examples of business protocols are those realized in cross-
11 organizational business processes, such as for negotiation, sales and purchase,
12 outsourcing of various business functions, delivery, repair, and so on. Traditionally,
13 business protocols have been modeled in purely operational terms
14 such as through state transition diagrams or message sequence charts that
15 describe ordering and occurrence constraints on the messages exchanged, but
16 not the meanings of such messages.

17 Commitments provide a natural basis for capturing the meanings of the
18 messages. In this manner, they provide a standard of correctness. A participant
19 in a business protocol complies with the protocol if it ensures that if any
20 commitment (of which it is the debtor) is detached, then it is also discharged
21 (that is, not violated or canceled—neglecting the distinction between them).
22 Having such a declarative basis for correctness not only simplifies the modeling
23 of the interactions being designed or analyzed but also provides a basis for
24 flexible enactment that can be shown to be correct.

25 A typical use of commitments in business protocols involves introducing
26 the syntax for the messages under consideration along with a formalization of
27 the meanings of the messages expressed in terms of the commitments of the
28 participants and the domain or environmental propositions that have a bearing
29 on those commitments. For example, in a purchase order protocol, we might
30 introduce a message *offer* and define its meaning as involving the creation
31 of a commitment—with its sender being the debtor and its receiver being the
32 creditor. The commitment would specify that the sender would provide the
33 goods to the receiver were the receiver to accept the terms. Likewise, we
34 might introduce a message *accept* through which the recipient of the *offer*
35 would take up the given offer. Based on these meanings, we would be able
36 to determine if an enactment of the protocol was sound. Even a simple and
37 obvious specification of correctness proves effective: this states that an agent
38 complies with a protocol if no enactment of the protocol ends with the agent
39 as the debtor of a detached but not discharged commitment.

40 The natural connection between commitments and correct enactments nat-
41 urally leads to ways of operationalizing them. Each commitment provides a
42 basis for judging the compliance of its debtor. The commitments of interest
43 taken together provide a *public* or neutral perspective on the correctness of an

1 interaction as a whole. Further, the idea of using both commitments that refer
2 to the antecedents and consequents of other commitments and commitments
3 that refer to the creation and manipulation of other commitments provides a
4 powerful basis for capturing a network of social relationships at a high level.
5 An agent can thus reason about the commitments of interest to it, especially
6 those where it is the debtor or creditor, and decide how to interact with the
7 other agents participating in the current business protocol. Although the agent
8 may act as it pleases, the commitments themselves impose constraints in terms
9 of what actions are compliant. In this sense, the specification of commitments
10 leads to the notion of a commitment machine [10, 51, 53, 42].

11 When we apply commitments as a basis for the semantics of the communi-
12 cations among agents, they yield a basis that is formal, declarative, verifiable,
13 and meaningful [41]. Interestingly, commitments also lend themselves to op-
14 erationalization in a more traditional manner. This is the idea of compiling
15 a commitment machine into a traditional representation such as a finite state
16 machine over finite [10, 51, 53] or infinite [42] computations. Such compilation
17 removes the opportunities for flexibility that an explicit commitment repre-
18 sentation supports. However, a finite state machine can be executed by agents
19 who are not equipped with an ability to reason logically. Moreover, such a
20 mechanically produced finite state machine can often be more complete in its
21 coverage of important scenarios than a hand-generated one—and consequently
22 be too large and unwieldy for a human designer to specify by hand.

23 3.2 Commitments for Communication Languages

24 The above idea involving protocols can also potentially be applied as a basis
25 for the meanings of the primitives in agent communication languages (ACLs).
26 ACL primitives have traditionally been given semantics based on the beliefs
27 and intentions of the communicating parties. Instead, a commitment-based
28 semantics could naturally express the social relationships between the com-
29 municating parties. In essence, one would take the idea of commitments for
30 individual communication protocols and apply that idea to the modeling of
31 general-purpose communication primitives. The idea is not implausible in it-
32 self. It is indeed possible to define the meanings of communication primitives.
33 In spirit, this is not different from the meanings of the messages in the busi-
34 ness protocols. However, the particular formulations in this setting suggest
35 ways to capture richer subtleties of meaning than may be necessary in a typ-
36 ical business setting. In particular, I have suggested [41] that meanings can
37 be captured via a trio of specifications that, following Jürgen Habermas [24],
38 reflect objective, subjective, and practical meanings. These types of meaning
39 can be expressed in terms of commitments regarding, respectively, the rele-
40 vant aspects of objective, subjective (cognitive), and practical (subjective and
41 institutional, with an emphasis on the latter) reality.

42 For example, we might define an *informative* message type as one creating
43 a dialectical commitment with its sender as debtor, its receiver as creditor, its

1 antecedent as true, and its consequent as asserting the truth of the proposition
2 specified as the content of the message. In the above terms, this would be the
3 objective meaning [41]. Likewise, a *commissive* message type would create a
4 practical commitment. And, similarly for the rest. I should note in passing that
5 the idea of a general-purpose ACL itself is suspect (see my recent manifesto,
6 as included in [15], for a discussion of this point). In any case, we can view
7 the definitions of the primitives as useful patterns, which might be specialized
8 and applied to the communicative acts needed for particular protocols.

9 *3.3 Commitments and Conventions*

10 A deeper benefit of commitments is in their relationship to conventions. Two
11 levels of abstraction are worth distinguishing in formalizing even the simplest
12 interactions. First, a quote means that there is a commitment from the merchant
13 to sell the specified goods at the specified price. Second, the fact that the
14 quote means the above is a matter of convention in the chosen domain of
15 commerce, and therefore both the merchant and the customer commit to that
16 meaning. Specifically, the meanings of any communications must be based
17 on the conventions at play in the given social setting. It is thus highly natural
18 that we understand conventions as a basis for interactions among autonomous
19 parties.

20 In several cases, the applicable conventions would be determined based
21 on longstanding tradition in a domain; in other cases, they may be explicitly
22 negotiated. For example, in the financial capital markets, a price quote for a
23 stock (sent by a broker to a trader) is interpreted as being merely informative
24 of the last known price at which that stock was traded. In typical commerce,
25 however, a price quote (sent by a merchant to a customer) can be interpreted
26 as an offer to sell at the specified price. In the latter case, the longevity of the
27 offer can vary: for a business-to-business supply price quote, the offer may
28 be valid for 30 days whereas for an airline to consumer ticket price, the offer
29 may be valid for a minute. The longevity of the offer too is often a matter
30 of convention. The importance of conventions to meaning and interoperation
31 among autonomous parties is thus quite obvious.

32 What is interesting for us is that the conventions that arise in a given setting
33 can be expressed as commitments. Specifically, each of the parties involved
34 (or sufficiently many of them) would commit dialectically to the existence
35 of the convention. Dialectical commitments, as are involved in this case, are
36 different from the practical commitments involved in formalizing the messages
37 in typical business protocols. However, each party may additionally practically
38 commit to acting according to the conventions. Arguably, something prevails
39 as a convention in a community only if the participants dialectically commit
40 to it and practically to acting according to it.

41 The general notion of conventions and especially as related to agent com-
42 munications [28], however, merits study in its own right. The interesting obser-
43 vation from the standpoint of commitments is that a convention corresponds

1 to an aggregation of dialectical commitments. The commitments can be struc-
2 tured using the context as explained above. Thus the participants in a commu-
3 nity where a convention prevails can dialectically commit to the convention.
4 Each participant would be a debtor and each other participant would be a
5 creditor. Alternatively, the creditor could be the context and thus stand for the
6 community as a whole.

7 **4 Concomitant Confusions**

8 In the worlds of artificial intelligence and software engineering abstractions,
9 commitments are the new kid on the block. A common prejudice in these
10 traditional disciplines that finds its way into multiagent systems is to formulate
11 the design problem as one for a complete unitary system, even when such a
12 system is to serve the needs of multiple stakeholders. Hence, all too often,
13 researchers and practitioners approach the design of a multiagent system not
14 only as consisting of cooperative (and sincere) agents, but also as one where
15 they will themselves provide all of the agents.

16 In contrast, commitments are most germane and offer their greatest value
17 in settings where capturing the meanings of the interactions being designed
18 is relevant. We would leave the design and construction of the agents to their
19 implementers even though in some cases we might ourselves take on the
20 implementation task. Further, we would leave the operation of the agents to the
21 agents and their users. That is, commitments can apply in traditional settings
22 where all agents may be designed by one party, and can help specify cleaner
23 architectures. But they are not confined to such settings, and the assumptions
24 needed for a unitary system do not apply in general to commitments.

25 One can imagine an engineer thinking “well, I am going to design a good
26 system of three agents; I am going to make sure the agents take on goals and
27 beliefs that are compatible with their commitments and adopt policies that
28 help them realize their commitments; and I am going to damn well make
29 sure the agents walk the straight and narrow, so I will prevent them from
30 violating their commitments.” Such thoughts may well be appropriate in a
31 single-perspective, cooperative, regimented system constructed by one engi-
32 neer from a set of agents. I would place the work of Minsky and Ungureanu
33 [33] into this category who are not focused on cognitive agents but on con-
34 ventional architectures, in which setting their approach is more reasonable.
35 However, such thinking unnecessarily limits the multiagent systems designs
36 that one comes up with. Therefore, although such thinking may be a useful
37 design pattern to help relate open architectures to traditional architectures,
38 when framed as a general constraint on commitments, it is misguided.

39 In simple terms, we can separate three scopes of effort or decision making:
40 (1) the modeler of an interaction defines interactions via their associated com-
41 mitments; (2) the agent designer implements an agent; and (3) the agent (and
42 its users) decide how to behave on the field. The multiagent system engineer
43 must specify the interaction precisely and relinquish control of the design and

1 operation of the endpoints of the interaction. Relinquishing control is a con-
2 sequence of dealing with open systems. Focusing on interactions is the only
3 plausible way of engineering a system where the engineer lacks control over
4 the endpoints.

5 *4.1 Commitments versus Goals*

6 A common view is that an agent who commits as debtor to bringing about
7 a condition in the world also adopts the same condition as a goal. (In some
8 accounts, the agents would adopt an intention, not just a goal, but let us
9 disregard the distinction between goals and intentions here.) A stronger variant
10 is when the goal applies to both the debtor and the creditor of a commitment.
11 This confusion is insidious because it relies upon a careless reading of the
12 literature: the confusion is nothing more than a confusion between the S-
13 commitments (our commitments here) and the P-commitments (traditional
14 commitments in AI planning) as explained in Section 2.

15 Commitments and goals are fundamentally different kinds of creatures.
16 A commitment is a public or observable relationship between two parties
17 whereas a goal is a single-agent representation. An agent's commitments are
18 generally known to others because of the conventions in play in the given set-
19 ting. An agent's goals are never inherently known to another agent, although
20 another agent might reason about them based on assumptions of rational-
21 ity or based on explicit revelation by the first agent, provided appropriate
22 conventions apply to the presumed revelation.

23 It is true that in general a cooperative debtor that created a commitment
24 would simultaneously adopt the corresponding goal. However, an agent
25 may not adopt the corresponding goal, potentially risking failing with its
26 commitment—and thus risking harm to its reputation and risking additional
27 sanctions of penalties and censure. Conversely, an agent may hold a goal and
28 not have committed to any other party for it. Such a goal might well be a
29 highly important goal for the agent—after all, a goal would relate to the agent's
30 preferences, and not necessarily to something the agent would reveal to others.

31 As a telling example, consider the common situation where an airline op-
32 erating a 100-seat airplane books 120 passengers on it. Clearly, the airline is
33 committed to each ticketed passenger, but equally clearly the airline could not
34 have a goal to board each passenger on to the airplane. The airline simply has a
35 clever internal strategy to maximize profit where it knowingly enters into com-
36 mitments that it might not be able to discharge. If 20 passengers miss the flight,
37 the airline goes scot-free; if more than 100 show up, the airline compensates
38 those it does not take on board, but it still comes out ahead on average.

39 Misalignments between commitments and goals are not the same as decep-
40 tion. In the above example, the airline has no intention of misleading its clients.
41 In fact, the airline may strongly believe—based on the evidence at hand—that
42 no more than 100 passengers will show up and thus none of its commitments
43 would be violated. However, it is fair to say (as a reviewer suggests) that a

1 commitment that is supported by its debtor's goal would be likelier to be
2 effective provided the debtor is sufficiently competent.

3 *4.2 Commitments versus Beliefs*

4 It is not uncommon to conflate commitments with beliefs. The motivation
5 seems to be that an agent would represent its commitments and thus believe
6 them to exist. But such an argument would hold for just about any represen-
7 tation.

8 In some cases, there is a more subtle confusion between commitments of
9 the dialectical flavor and beliefs. Notice that even dialectical commitments
10 are commitments, meaning that they reflect their debtor staking a claim or
11 accepting a claim as a putative fact such as, for the sake of discussion, during
12 an argument [32]. The debtor may not in fact believe what it commits to.
13 Conversely, the debtor may have numerous beliefs it keeps private and never
14 commits to holding to another agent. Any such commitment binds the debtor
15 to a certain pressure to interact in a certain way, and there is generally little
16 reason to expose all beliefs as commitments.

17 *4.3 Commitments versus Mutual Beliefs*

18 A more insidious confusion arises with respect to mutual beliefs. As Section 2
19 explains, the underlying idea behind using the mutual beliefs (and equally
20 intentions) was to introduce a level of mutuality while continuing to use the
21 mental concepts.

22 The first problem with this view is that it is wrong. Commitments are not
23 mutual beliefs. A commitment is a unidirectional relationship. For example,
24 if Bianca offers to sell a camera to Alessia, the commitment holds whether
25 or not Bianca believes it or Alessia believes it. As in the airline over-booking
26 example, Bianca may simply have made the offer to try to prevent Alessia from
27 taking up another offer. And Alessia might be on to Bianca: that is, she might
28 not believe that Bianca believes she would supply the camera. However, the
29 commitment exists. Alessia may in fact file a complaint against Bianca. Alessia
30 would not be able to file a complaint if the commitment was defined as the
31 mutual belief.

32 The second problem is that mutual beliefs are extremely fragile. Let us say
33 Alessia believes that Bianca believes that Alessia believes . . . that Bianca will
34 be shipping a camera to Alessia. If Alessia believes that, at the hundredth level
35 of the nesting, Bianca might not believe Alessia expects the camera any more,
36 that would dissolve the mutual belief. However, in real life the commitment
37 does not go away in such a case. Bianca is not off the hook based on a failure
38 of a belief by Alessia and certainly not for imagining that Alessia may have
39 lapsed in her belief.

1 The third problem is that, again at variance from real-life interactions, al-
2 though commitments arise in all manner of distributed settings, mutual be-
3 liefs generally cannot be constructed. Under asynchronous communication,
4 the only mutual beliefs in a system are the invariants of the system, that is,
5 propositions that were true from the start [8]. Indeed, the artifact of mutual
6 beliefs (along with the similar artifact of common knowledge) is used in dis-
7 tributed computing primarily to prove *impossibility* results [8, 25]. Because
8 mutual beliefs cannot be engendered through message exchange in general
9 asynchronous settings, a problem that requires mutual belief is unsolvable.

10 Clearly, the AI researchers have understood the problem in terms of live hu-
11 man communication, which is inherently synchronous. In multiagent settings,
12 they address the challenges of asynchrony by fiat. Specifically, they assume
13 that a single message by one party to another, without any need for an ac-
14 knowledgment, would achieve mutual belief. The idea it seems was that there
15 was a central belief store and any assertions inserted into it reflected the beliefs
16 and further even the mutual beliefs of everyone in the system. However, AI
17 researchers by and large hide this key assumption in the implementations of
18 their systems and never mention it in their theoretical descriptions.

19 One could treat the above assumption (of a single message exchange being
20 sufficient) as a standard operating procedure, as mentioned in Section 2, in
21 a particular setting. But that only means we are seeking to characterize com-
22 mitments a certain way. So why not be honest and model the commitments
23 directly? About the only reason not to do so is if one has locked on to the
24 mentalist ideology.

25 I should explain that the point is more general than merely one of physical
26 transmittal of information, as it is in the traditional distributed computing liter-
27 ature. The deeper and more crucial point is of the necessity of simultaneously
28 sustaining multiple perspectives. In other words, what is most problematic
29 is not so much the physically central nature of the belief store where mutual
30 beliefs might exist, but its conceptually central nature, indicating that we had
31 magically consolidated the perspectives of multiple autonomous, heteroge-
32 neous parties into a correct unitary perspective.

33 **4.4 Commitments versus Obligations**

34 Obligations are an important notion studied since ancient times. A traditional
35 obligation applies on an agent, roughly corresponding to the debtor of a com-
36 mitment. What distinguishes a traditional obligation from the cognitive con-
37 cepts of beliefs and goals is that it is inherently externally focused: an obligation
38 can be met or not and the consequences occur beyond the minds of the agents
39 involved. A more interesting kind of obligation is directed: here an agent is
40 obliged to another agent [26]. The second agent corresponds to the creditor of
41 a commitment.

42 Because directed obligations are clearly interagent in their orientation, they
43 are a more natural match for multiagent systems than are traditional obli-

1 gations. The similarities between directed obligations and commitments are
2 striking. But can we treat commitments as being identical to obligations? A
3 commitment when it is active corresponds to a directed obligation.

4 However, commitments and obligations have important points of distinc-
5 tion. First, a commitment can be manipulated, in particular, delegated, as-
6 signed, or released. Second, a commitment carries with it an organizational
7 context, as explained above. Third, obligations carry a moral connotation that
8 commitments lack. Fourth, a commitment reflects the inherent autonomy of
9 the participants in an interaction. Thus an agent would become a debtor of
10 a commitment based on the agent's own communications: either by directly
11 saying something or having another agent communicate something in con-
12 junction with a prior communication of the debtor. That is, there is a causal
13 path from the establishment of a commitment to prior communications by
14 the debtor of that commitment. Obligations by contrast can be designed in or
15 inserted by fiat.

16 Frank Dignum observes (in a private comment) that the autonomous nature
17 of commitments raises a creditor's expectation that the debtor's goals and be-
18 liefs are aligned with the commitment, and hence it should be discharged. This
19 point applies to cooperative debtors and may be a basis for the conventional
20 interpretation of communications in general.

21 *4.5 Commitments versus Policies*

22 A commitment, especially in its conditional form, looks like a rule for pro-
23 cessing, and in this sense resembles a policy. For example, an engineer might
24 take the view that an offer from a merchant to a customer expresses the mer-
25 chant's policy that if the customer pays the specified amount to the merchant
26 the merchant will send a cello string of a specified type to the customer.

27 Treating a commitment as a policy in this sense reflects the same confusion as
28 with goals and beliefs, namely, that the external, interactive, observable nature
29 of commitments is conflated with the internal, behavioral, private nature of
30 another abstraction. A policy is how an agent may decide to act upon—or decide
31 not to act upon—a commitment. If the merchant has a straightforward policy
32 for acting on all its commitments, then so much the better.

33 However, note that in general, a commitment would be necessarily in-
34 complete with respect to the behavior needed to discharge it, and thus the
35 policies associated with a commitment may need to specify aspects that the
36 commitment does not mention. In our example, the merchant would have
37 committed simply to supplying, say, a Larsen cello D string for payment. The
38 merchant would need additional policies to determine how exactly to supply
39 the D string. Should the merchant supply the instance of the D string that is
40 the oldest in its inventory? Or, the newest? Or, one that happens to be the most
41 convenient based on other tasks the merchant is performing, for example, sup-
42 ply from the top rack if the ladder is up there anyway, else supply from the
43 bottom rack? Maybe the merchant will do well to supply a carefully checked

1 instance of the string to a repeat customer and supply it with extra robust
2 packaging for a customer overseas or for a customer who has a high standing
3 in the user community and can influence other prospective customers.

4 These are all legitimate policies, but it would be inappropriate to tie them
5 into the commitments. Indeed, were we to attempt to specify commitments at
6 the level of such policies, we would face important challenges and produce
7 a poorer quality model as a result. The challenges would be first coming
8 up with detailed specifications and second, importantly, finding a way to
9 determine if a party is complying with the commitment—for example, how
10 will we ever know if the merchant sold the oldest item in its inventory? Or,
11 one from the top rack? The resulting model would be of poor quality because
12 it would tightly couple the parties involved in the interaction. For example, a
13 merchant who committed to supplying the oldest item (and did so honestly)
14 would be compelled to maintain information about the ages of the items in its
15 inventory and to set up its internal business processes to search for items in
16 their order of age. It would not be able to take advantage of any improvements
17 in internal business processes as might arise later. Equally importantly, a new
18 merchant who wished to join the interactions specified by such a commitment
19 would not be able to participate without developing such otherwise irrelevant
20 components of its information systems.

21 There is another notion of policies, however, which does make sense when
22 related to commitments. This is the idea of a social policy, which captures the
23 rules of encounter in a society. I have occasionally used the term “social policy”
24 in this sense, but I now think it is better to refer to such as *norms* and to reserve
25 the word “policies” for the policies of an agent or organization that reflect its
26 decision making.

27 *4.6 Commitments versus Regimentation*

28 I encounter this problem a lot in discussions with conventional software engi-
29 neers. They are accustomed to capturing requirements for, modeling, design-
30 ing, and implementing software systems in which there is a single locus of
31 autonomy. The system in question may be distributed but it is conceptually
32 unitary and involves the perspective of a single party. You can identify such
33 a mindset where the engineer talks of “the system” as an entity that will in-
34 teract with “the user”—the goal of the engineer is to create a software design
35 artifact from which one can develop a set of software modules that will meet
36 the elicited requirements as the system interacts with its users.

37 In such a case, when the engineer begins reluctantly to think of social inter-
38 actions and commitments among the parties involved, the engineer’s mindset
39 remains to try to force the modules to behave in the “correct” manner. The en-
40 gineer attempts to capture such behaviors via commitments. In other words,
41 the engineer retains the single-party perspective and, without absorbing the
42 idea of any social interaction among notionally autonomous parties, merely
43 treats commitments as a clever-sounding representational framework.

1 The engineer's challenge is to force the modules to adopt certain commit-
2 ments and to act on their commitments in exactly the chosen "correct" way.
3 All too often, such designs emerge from when a traditionally minded soft-
4 ware developer reverse-engineers an existing process into the representation
5 of commitments—adopting and incorporating every ad hoc quirk of the original
6 model into the commitment-based model.

7 Following Jones and Sergot [29], I term such a viewpoint *regimentation*. In
8 general, the use of regimentation obviates the need for modeling commitments.
9 However, for engineers new to commitments, it might be a useful intermediate
10 step provided they recognize it as such and proceed to develop an interaction-
11 oriented model.

12 **4.7 Commitments and Compliance**

13 When computer scientists and business (process) modelers first encounter
14 commitments, they immediately ask about compliance: how can we guarantee
15 that an agent would comply with its commitments, or at least not wantonly
16 violate or cancel them? For a novice, this question is reasonable. But upon
17 reflection, we can see this question is misguided and unfair because it hides
18 some crucial presuppositions and confusions. Underlying this question is the
19 misguided assumption that if one simply fails to model—or even acknowledge
20 the existence of—an agent's commitments, the agent would behave perfectly.

21 A strange variation on this theme is that if we were to model communi-
22 cations among agents in terms of commitments we would have created legal
23 liabilities that didn't exist before. No, seriously, I am not making this up. The
24 idea is that if Bianca sends Alessia a message with an offer for a camera, for
25 example, using English or XML, it is just fine and legally safe. But if we only so
26 much as realize that the offer is a commitment to provide the specified camera,
27 Bianca would become liable in ways that she wasn't when we didn't model
28 her English or XML message as conveying a commitment. Perhaps the people
29 who come up with the above variants imagine that obfuscation of meaning
30 is a legal defense. I claim, instead, that for a business or other interaction to
31 be successful, the parties involved must share an understanding of the terms
32 involved. In general, even lawyers would prefer greater clarity as a way to
33 define each party's expectations of the others.

34 Modeling commitments does not cause agents to (potentially) behave in an
35 undesirable manner. Indeed, modeling commitments helps potentially address
36 the challenge of ensuring compliance. By treating commitments explicitly, we
37 (1) obtain a crisp, yet not operational, statement of compliance; (2) formulate
38 the notion of *transparent* protocols in which compliance determination is possi-
39 ble; and (3) open the way for designing agents using beliefs and goals who
40 will be compliant with their protocols. Monitoring and compliance relate nat-
41 urally to themes such as formalizing (1) organizations and governance [50, 21],
42 for example, penalizing malfeasant agents in a community, and (2) bases for
43 relating commitments and economic models of rationality [16].

1 **4.8 Terminological Confusions**

2 It is worth highlighting here some confusions that arise alarmingly frequently,
3 though usually among people who are unfamiliar with the commitments lit-
4 erature. At the root of these confusions are lexical mismatches, wherein the
5 reader misinterprets a technical term, even though the terms under consider-
6 ation are well-defined in the commitments literature.

7 Commitments are psychological. The *raison d'être* for commitments is to
8 avoid the shortcomings of psychological commitments, but that doesn't
9 stop some people from inadvertently going back to square one.

10 Social means going to a bar. We use the term social to distinguish from psy-
11 chological, not that commitments are only about cultural conditions or for
12 after-hours socialization. The most common application of commitments
13 today is in modeling business organizations and interactions, though there
14 is no reason to preclude other settings even personal relationships.

15 Private means shared. Private refers to the internals of an agent and public
16 to what is shared or observable. If one agent commits to another, that means
17 we have created a social object involving at least two agents. Even if the
18 two agents keep the commitment confidential, never disclosing it to a third
19 party, the fact that it involves more than one agent makes it public, as we
20 define the term.

21 Debts are exclusively financial. We simply use debtor and creditor to indi-
22 cate the directionality of commitments. These terms are reminiscent of their
23 usage in the vernacular, but generalize over it. There is no restriction to
24 financial debt: the conditions involved could be arbitrary; indeed, even in
25 normal English, debts are not restricted to be just financial.

26 Organizational context means any element of the situation. But organizational
27 context is *not* just anything: in our technical meaning, it is an objective in-
28 stitutional construct treated on par with an agent.

29 Commitments are ontological commitments. In Quine's [34] terminology,
30 an ontological commitment describes the objects one entertains as exist-
31 ing. For example, if I say my grandfather owned a unicorn, that means I am
32 ontologically committed to the past existence of at least one unicorn (and of
33 my grandfather, and of the two existing contemporaneously). Ontological
34 commitments resemble presuppositions underlying utterances that a per-
35 son makes whereas commitments for us are about actions or staked claims.
36 One could formulate a dialectical commitment for the existence of anything
37 that its debtor makes an ontological commitment to.

38 **5 Debate with Professor Castelfranchi**

39 Let me now turn to the most interesting part of this article, which is to highlight
40 some of the points of difference between Professor Castelfranchi's views and
41 mine regarding commitments.

1 Let me begin with a point, which I think is not controversial, though po-
2 tentially sounding like it might be. At the root of it is the emphasis I place
3 on the importance of observable interactions among agents (including low-
4 level behaviors), which contrasts with Professor Castelfranchi's emphasis on
5 the cognitive representations of the agents. I suspect unclarity on my part led
6 Professor Castelfranchi to criticize my approach as resembling a behaviorist
7 approach.

8 A behaviorist stance would be reduced to entertaining nothing beyond
9 (what the designer or analyst imagines are) the objective atoms of behavior.
10 In general, the difficulty in identifying such objective atoms is indeed one of
11 the challenges that uproots behaviorism. The acute need for imagining what
12 is ostensibly objective is one of the shortcomings of behaviorism. However, I
13 do not see that commitments can be reduced merely to low-level behaviors.
14 Instead, here we are accommodating a rich social reality: we have postulated
15 agents who create and function in social institutions, who entertain abstract
16 high-level relationships such as those expressed via commitments, and who
17 not only communicate at the level of exchanging bits of information but also
18 communicate in suitable institutional terms.

19 Professor Castelfranchi and I thus agree that the study of commitments is
20 not and should not be treated as a behaviorist project. Instead, our collective
21 effort in multiagent systems may be thought of as a realist project in that we
22 treat common-sense social constructs such as commitments as real entities.

23 *5.1 Commitments and Autonomy*

24 Broadly speaking, the multiagent systems field is primarily concerned with
25 understanding the interactions of agents. At a basic level the autonomy of
26 the agents is key. Of course, fully autonomous agents would be useless if not
27 harmful—clearly, what we need to understand is the interdependence of the
28 agents. That is exactly where commitments come in. Each commitment cap-
29 tures one element of a social relationship between two parties. When we put
30 these elements together, we obtain the network of relationships that character-
31 izes a multiagent system. I expect that Professor Castelfranchi and I agree on
32 the above in broad terms.

33 Where I suspect we disagree is in the relative importance we accord the
34 intuitions of autonomy and interdependence. As I see it, an agent must be able
35 to enter into and exit its commitments at will whereas Professor Castelfranchi
36 sees the process as more constrained. These distinctions become more apparent
37 when we consider the creation or cancellation of a commitment.

38 **5.1.1 Accepting a Commitment**

39 Professor Castelfranchi sees a commitment in a positive light whereas I see
40 it as a general notion in a neutral light. Also, my interest is to maximize the

1 flexibility of the interactions and the autonomy of the participants. As a result,
2 I would consider a commitment to be created if its debtor says so. In this sense,
3 the creation of a commitment is a declarative or performative communication
4 and is within the control of the agent initiating that communication, given
5 the appropriate circumstances and conventions. In contrast, Professor Castel-
6 franchi would like to see the creditor of a commitment explicitly accept the
7 commitment before it comes into being.

8 A downside to Professor Castelfranchi's approach is that it couples the two
9 agents unnecessarily. It also differs from common uses of commitments. For
10 example, a merchant can make an offer to a customer merely by saying so.
11 The customer may sit silently for a while (up to the time period of the offer)
12 and then attempt to make a purchase based on that offer. That is, the customer
13 doesn't separately accept the offer and then exercise it; the customer simply
14 exercises the offer directly. The offer is valid all along. If we were to require that
15 the offer be accepted before it comes into existence, that would seem to require
16 that a message exchange has to complete before the offer begins to exist.

17 Professor Castelfranchi is concerned that if we do not include an explicit
18 acceptance, an agent may in essence use a commitment to make a threat, for
19 example, by committing to harm the creditor. In Professor Castelfranchi's ap-
20 proach, the creditor would refuse such a commitment and thus never let it
21 be formed. Notice, however, a malicious (prospective) debtor could harm the
22 creditor nevertheless. If the commitment happens to be undesirable for the
23 creditor, it could (i) resist it in other ways, perhaps by making a threat of its
24 own; (ii) ignore the commitment and not demand that the debtor discharge it;
25 (iii) assume it arose due to some underlying confusion due to miscommunica-
26 tion with the debtor, and explicitly release the debtor from that commitment.
27 Each of these sample approaches has the advantage of not creating avoidable
28 coupling between the debtor and the creditor.

29 Also, the apparent undesirable-to-the-creditor orientation of the content
30 of a commitment cannot always be avoided. For example, an organization's
31 president Alessia may have committed to all its members that she would
32 punish the treasurer were the treasurer to embezzle any funds. A member,
33 Bob, may accept the commitment at a meeting along with the other members
34 of the organization. Now later if Bob becomes the treasurer, he would be the
35 creditor of a commitment from the president that might potentially penalize
36 him, if it is activated at all.

37 An alternative view is that the above notion of acceptance ought to be
38 considered as being explicit *or* implicit. Thus silence in our example above
39 can be treated implicit consent. This view, however, misses two important
40 points. The first point is that it contravenes the agents' autonomy, as explained
41 above. The second point we can explain as follows. The deeper purpose of
42 talking about commitments is to help us understand the social state of an
43 interaction. If we decide that a commitment is created only upon acceptance
44 by the prospective creditor that means we can provide no clear meaning for the
45 intermediate state wherein the debtor has "committed" but not quite because
46 the creditor has not confirmed yet. If we allow implicit acceptance, then we
47 have no viable basis for distinguishing between the commitment and its half-

1 baked stage. That half-baked commitment is not nothing because the debtor is
 2 on the line if the creditor accepts it. I claim that if the associated intermediate
 3 social state were to be formalized properly, the semantics that results from the
 4 acceptance-based approach would be close to that of the one-sided formulation
 5 that I advocate.

6 Consider the following example, which came up in a discussion with Neil
 7 Yorke-Smith. How might one model the following? Alessia proposes to Bob
 8 that they exchange goods for payment tomorrow, but today Alessia would like
 9 to know whether Bob accepts or not.

10 A simple formulation is $C(\text{Alessia}, \text{Bob}, C(\text{Bob}, \text{Alessia}, \text{goods}, \text{pay}), \text{goods})$,
 11 indicating that Alessia tells Bob “if you commit to pay on receipt, I will send you
 12 the goods.” It’s Alessia’s decision to trust Bob. If Bob does commit, Alessia must
 13 send the goods or violate her (now detached) commitment. If Alessia sends
 14 the goods after Bob’s acceptance, Bob must pay or violate his (now detached)
 15 commitment. This formulation shows how we can make the acceptance of
 16 a commitment explicit if and when we need it to model some scenario, but
 17 do not need to insist upon acceptance in other cases. We can think of the
 18 above formulation as interpolating two one-sided commitments Alessia to
 19 Bob: one conditional on payment, $C(\text{Alessia}, \text{Bob}, \text{pay}, \text{goods})$ and the other
 20 unconditional $C(\text{Alessia}, \text{Bob}, \text{true}, \text{goods})$. In contrast, the acceptance-based
 21 representation makes it impossible to express the one-side commitments; tends
 22 to be applied wrongly wherein one agent commits another, thereby violating
 23 the latter’s autonomy; and, leaves as undefined the social state wherein Alessia
 24 has made an offer but Bob hasn’t responded.

25 5.1.2 Accepting a Cancellation

26 In much the same spirit, I propose that an agent can cancel its commitment at
 27 will. Like creation, a cancellation is a declarative that the debtor can perform.
 28 Likewise, a creditor can perform the release of a commitment at will. In the case
 29 of the cancellation, the outcome might not be one that the creditor desires or
 30 would willingly accept; further, the outcome might be one that we as designers
 31 might not condone in our agents. However, if that were to be the case, the
 32 creditor should have made sure (or we, the designers, should have made
 33 sure) that there will be repercussions on the debtor for having performed an
 34 inappropriate cancellation.

35 One might think that these repercussions signal the unacceptability of the
 36 cancellation and, therefore, that cancellations should only be allowed when the
 37 creditor accepts. I won’t repeat the points made above in connection with cre-
 38 ating a commitment, which apply here too. However, an additional point rel-
 39 evant to cancellation is that in a multiagent system (consisting of autonomous
 40 agents), we can rely on *regulation* but not on *regimentation* [29]. Regulation
 41 is about controlling behavior through normative means whereas regimenta-
 42 tion is simply about preventing bad behavior [1]. Regulation is suited to in-
 43 teractions among autonomous agents. In contrast, regimentation—which here

1 corresponds to preventing cancellation by explicit acceptance—contravenes au-
2 tonomy.

3 Even in the original formulation of commitments [38], the notion of the
4 (organizational) context of a commitment served to accommodate such cases.
5 Specifically, if the cancellation of a commitment arises because of true and
6 reasonable exceptions, the context may impose no penalty upon the debtor;
7 in other cases it might. For example, let's say a merchant has committed to
8 providing some goods to a customer. If the merchant cancels the commitment
9 to do so because of a tsunami that destroyed the manufacturing plant and
10 refunds the customer's payment, the cancellation appears not unfair whereas
11 if the merchant cancels because the merchant can now demand a higher price,
12 the cancellation does sound egregious. Let us say the (organizational) context
13 here is the electronic marketplace, for example, eBay. In the first case, the
14 context may declare the cancellation legitimate; in the second case, not so.
15 In the second case, the context may penalize the merchant, for example, by
16 revoking his credentials in the marketplace or pursuing fraud charges in the
17 court system.

18 If the organizational context can ensure such coherent outcomes, then we
19 can think of the context (and the concomitant family of interactions) as being
20 well-designed (notice we make no claims about the internals of the agents
21 themselves). If the context is not well-designed, then either we as designers
22 made a mistake or the agent (customer) made a mistake in joining such a
23 context, dealing with an untrustworthy merchant, and foolishly counting on
24 him to discharge his commitments.

25 *5.2 Commitments and Cognition*

26 Another of the points where I continue to have a disagreement with Professor
27 Castelfranchi is in the function and importance of cognitive representations
28 in connection with commitments. We agree, of course, on the basic idea that
29 an agent's behavior is of central importance in judging whether or not it dis-
30 charges its commitments. And, I expect we agree not only on the essential
31 relevance of commitments to the social life of an agent, including its relation-
32 ships with other agents, but also on the importance of cognition.

33 Professor Castelfranchi, however, assigns a far stronger function to the cog-
34 nitive representations of an agent than I do. To him, having a commitment
35 is strongly based on the associated patterns of beliefs, goals, and intentions.
36 For me, in contrast, a commitment is a social entity, which takes its existence
37 from the public sphere. An intelligent agent would undoubtedly represent
38 and reason about its commitments, and its commitments would undoubtedly
39 affect and be affected by its goals and intentions. However, to my thinking, a
40 commitment at its core remains purely social. In this regard, a commitment
41 is no more and no less of an abstract object than any cognitive attitude or
42 any mathematical object for that matter—that is, a commitment can exist in the
43 public sphere just as legitimately as in the mind of an agent.

1 Although I recognize the benefits and importance of the cognitive represen-
2 tations in modeling and implementing agents, I consider such representations
3 to be internal to an agent and reflective of its internal architecture and con-
4 struction. In contrast, I understand commitments as having normative force
5 whereby they can provide a potentially independent basis for judging the fe-
6 licity and correctness of the actions of agents. When we define commitments
7 in such a public and observable manner, they can become a key ingredient
8 in understanding the institutional nature of communications and indeed of
9 understanding institutions themselves.

10 As an example, consider a friend of mine who promises to help by giving
11 me a ride to the airport. My friend would have done so by using the prevailing
12 vernacular of our social institutions to create a promise. Let us say the ap-
13 pointed hour comes and goes, but my friend does not materialize. Thus he has
14 violated his commitment. For the sake of this example, let us further stipulate
15 both that I trust my friend in such matters and that he is highly trustwor-
16 thy in fact and would not have deceived anyone. Clearly, he forgot or found
17 himself in a personal emergency. But we would still state that he violated his
18 commitment, albeit inadvertently or in exonerating circumstances.

19 We should be able to pass the judgment of the commitment being violated
20 based on what we observe, namely, the failure of the commitment. However,
21 if the definition of commitments were to be intertwined with questions of
22 beliefs and goals, it would be difficult for us to pass even such elementary
23 judgments. Further, the definition would lose the benefit of modularity by
24 combining the social and the cognitive representations. Additionally, it would
25 create a situation where we would not be able to determine if a commitment
26 existed without being able to assess what the beliefs and intentions of the
27 parties involved were, and it is well-known that such ascriptions cannot be
28 defended in multiagent settings where the agents are not homogeneous and
29 their internal states not public [39].

30 I claim that such judgments provide the basis of the normative strength that
31 commitments carry. We might conduct any amount of elaborate post mortem
32 analyses involving the beliefs and goals of the participants, but if we are
33 not clear about the objective fact in this matter, we lose not only a basis for
34 specifying an institutional basis for multiagent systems but also for conducting
35 any cognitive analyses with any grounding in truth.

36 6 Themes for the Future

37 6.1 *Commitments and Trust in Social Computing*

38 The increasing attention garnered by topics such as social computing tells us
39 that areas of long interest in the multiagent systems field [22] and especially
40 pursued by Professor Castelfranchi himself [5] are gaining currency. Today's
41 practice in social computing is weak indeed and consists of little more than
42 users sharing information on a social networking site or users performing

1 various assigned tasks in what is called crowdsourcing. It seems to me self-
2 evident that any kind of realistic social computing must rely upon the concepts
3 of commitments and trust.

4 The study of trust has been an important theme in Professor Castelfranchi's
5 body of research. Professor Castelfranchi and colleagues have developed a
6 semantically rich notion of trust [6, 7, 18] that incorporates both its social
7 and its cognitive aspects. Professor Castelfranchi's approach contrasts with
8 the majority of computer science works on trust, which tend to jump into
9 (typically, numerical) representations without first sorting out what the trust
10 as conceived stands for. Professor Castelfranchi relates trust to the plans of the
11 parties involved and their expectations with respect to each other. I find another
12 of previous works by Professor Castelfranchi and colleagues as especially
13 germane here. This is the notion of dependence [36], which Rino Falcone and
14 Professor Castelfranchi [17] have recently revived and related to trust.

15 It seems clear to me that these concepts suggest the strong relationship
16 between commitments and trust. In conceptual terms, we can think of com-
17 mitments and trust as duals of each other: a debtor commits to a creditor and
18 a trustee places trust in a trustee. The idea of commitments as expectations
19 in reverse originates in Amit Chopra's [11] dissertation. I have recently be-
20 gun to formalize trust in a manner that highlights the notion of dependence
21 and relates trust to commitments [44]. Not every commitment may have cor-
22 responding trust in the reverse direction. And, not every placement of trust
23 may be justified by a commitment in the reverse direction. The best outcomes
24 arise when trust and commitment go hand in hand. The existence of trust for
25 a commitment means that the commitment is not superfluous. The existence
26 of a commitment for trust means that the trust is not misplaced. Chopra and
27 colleagues [14] investigate the connection of trust with architecture. Exploring
28 the above themes further and especially modeling social action as it would
29 arise in future application settings of even moderate complexity would be
30 highly valuable.

31 *6.2 Commitments and Software Engineering*

32 Let me now talk about another important theme with regard to commitments.
33 This has to do with the use of commitments in modeling and realizing multi-
34 agent systems in diverse domains. In today's practice, software engineering is
35 mainly concerned with low-level abstractions that are close to implementation
36 details. Such abstractions are difficult to specify and even harder to establish
37 the validity of with respect to the needs of the stakeholders.

38 Commitments provide a nice alternative basis for specifying software sys-
39 tems. Work on applying commitments for software engineering has been going
40 on for years, since the earliest studies, and initially under the rubric of commit-
41 ment protocols. However, the more basic challenges of software engineering
42 when applied to interactions in multiagent systems are now beginning to be
43 understood and formulated in terms of commitments [9, 12, 13, 31, 47].

1 Although the above approaches are useful and promising, they are far from
2 adequate when it comes to the challenges of building systems of practical
3 complexity. I foresee the enhancement of the techniques in terms of clearer
4 specification languages based on commitments, more extensive middleware
5 that supports implementation using abstractions similar to commitments, and
6 the development of tools and technologies to validate and verify commitment-
7 based designs.

8 In this light, I further think than commitments can inform an expanded
9 notion of norms. Unlike a lot of traditional work, wherein norms are treated
10 as amorphous descriptions of good or normative behavior, I propose that we
11 study norms that like commitments are directed, conditional, contextual, and
12 manipulable. Such norms can help precisely capture normative conditions in
13 a manner where it is clear who is responsible for their enforcement. The notion
14 of organizational context provides a basis for understanding the *governance*
15 of systems of autonomous parties [45], such as service engagements [50] and
16 virtual organizations [2, 48, 49].

17 7 Conclusions

18 I have taken this essay as an opportunity to lay out the main themes relating to
19 commitments. I imagine that Professor Castelfranchi and I largely agree with
20 each other on virtually all of the substantial themes regarding commitments.
21 I have highlighted some controversial points in the hope that they would be
22 interesting and useful, especially for those new to the field.

23 However, to summarize quickly, our points of agreement include the funda-
24 mental importance of understanding interaction in multiagent systems from
25 the social and institutional level as opposed to exclusively from the mechanical
26 or operational levels; the very conception of commitments as an elementary
27 social (as opposed to an exclusively mental relationship, as in AI); the dis-
28 tinctions and similarities between practical and dialectical commitments; the
29 value of commitments in understanding institutions and norms; the close re-
30 lationship between commitments on the one hand and dependence and trust
31 on the other.

32 Although the field of multiagent systems has made substantial progress
33 since its founding just decades ago, a lot of crucial theoretical and practical
34 problems remain unanswered and even unformulated. No one can predict with
35 any certainty where the field will grow. However, the emergence of networked
36 computing and its expansion into human business and social life suggests that
37 the future of multiagent systems—viewed as the academic field that studies the
38 interactions of social beings—is secure. That our field is now established and
39 has acquired a healthy respect for, if not yet universally a deep understanding
40 of, the social basis for interaction is due in no small part to the imagination
41 and intellect of one researcher and for these invaluable contributions I applaud
42 Cristiano Castelfranchi.

1 **Acknowledgements** I have benefited a lot over the years from discussions regarding com-
 2 mitments with a number of people, among them Cristiano himself and, alphabetically,
 3 Matthew Arrott, Alexander Artikis, Matteo Baldoni, Cristina Baroglio, Amit Chopra, Marco
 4 Colombetti, Nirmal Desai, Frank Dignum, Virginia Dignum, Rino Falcone, Nicoletta Fornara,
 5 Les Gasser, Scott Gerard, Paolo Giorgini, Kohei Honda, Michael Huhns, Andrew Jones, Mike
 6 Luck, Ashok Mallya, the late Abe Mamdani, Elisa Marengo, Simon Miles, John Mylopoulos,
 7 Viviana Patti, Jeremy Pitt, Pankaj Telang, Paolo Torroni, Yathi Udupi, Feng Wan, Michael
 8 Winikoff, Jie Xing, Pinar Yolum, and Neil Yorke-Smith. Comments from Michael Huhns,
 9 the Dignums, Scott Gerard, Pinar Yolum, and the anonymous reviewer have helped im-
 10 prove this article. I wouldn't presume, however, that any of the people named above agrees
 11 with anything I have claimed in this article. I would also like to thank the National Science
 12 Foundation for partial support under grant 0910868.

13 References

- 14 1. Artikis A, Sergot MJ, Pitt JV (2009) Specifying norm-governed computational societies.
 15 ACM Transactions on Computational Logic 10(1)
- 16 2. Brazier F, Dignum F, Dignum V, Huhns MN, Lessner T, Padget J, Quillinan T, Singh
 17 MP (2010) Governance of services: A natural function for agents. In: Proceedings of
 18 the 8th AAMAS Workshop on Service-Oriented Computing: Agents, Semantics, and
 19 Engineering (SOCASE), pp 8–22
- 20 3. Castelfranchi C (1993) Commitments: From individual intentions to groups and orga-
 21 nizations. In: Proceedings of the AAAI Workshop on AI and Theories of Groups and
 22 Organizations: Conceptual and Empirical Research
- 23 4. Castelfranchi C (1995) Commitments: From individual intentions to groups and orga-
 24 nizations. In: Proceedings of the International Conference on Multiagent Systems, pp
 25 41–48
- 26 5. Castelfranchi C (1998) Modelling social action for AI agents. *Artificial Intelligence* 103(1–
 27 2):157–182
- 28 6. Castelfranchi C, Falcone R (2010) *Trust Theory: A Socio-Cognitive and Computational*
 29 *Model*. Agent Technology, John Wiley & Sons, Chichester, UK
- 30 7. Castelfranchi C, Falcone R, Marzo F (2006) Being trusted in a social network: Trust as
 31 relational capital. In: *Trust Management: Proceedings of the iTrust Workshop*, Springer,
 32 Berlin, LNCS, vol 3986, pp 19–32
- 33 8. Chandy KM, Misra J (1986) How processes learn. *Distributed Computing* 1(1):40–52
- 34 9. Cheong C, Winikoff MP (2009) Hermes: Designing flexible and robust agent interactions.
 35 In: Dignum V (ed) *Handbook of Research on Multi-Agent Systems: Semantics and*
 36 *Dynamics of Organizational Models*, IGI Global, Hershey, PA, chap 5, pp 105–139
- 37 10. Chopra A, Singh MP (2004) Nonmonotonic commitment machines. In: Dignum F (ed)
 38 *Advances in Agent Communication: Proceedings of the 2003 AAMAS Workshop on*
 39 *Agent Communication Languages*, Springer, LNAI, vol 2922, pp 183–200
- 40 11. Chopra AK (2008) *Commitment alignment: Semantics, patterns, and decision proce-*
 41 *dures for distributed computing*. PhD thesis, Department of Computer Science, North
 42 Carolina State University
- 43 12. Chopra AK, Singh MP (2011) Specifying and applying commitment-based business
 44 patterns. In: Proceedings of the 10th International Conference on Autonomous Agents
 45 and MultiAgent Systems (AAMAS), IFAAMAS, Taipei, pp 475–482
- 46 13. Chopra AK, Dalpiaz F, Giorgini P, Mylopoulos J (2010) Modeling and reasoning about
 47 service-oriented applications via goals and commitments. In: Proceedings of the 22nd
 48 International Conference on Advanced Information Systems Engineering (CAiSE), pp
 49 417–421

- 1 14. Chopra AK, Paja E, Giorgini P (2011) Sociotechnical trust: An architectural approach.
2 In: Proceedings of the 30th International Conference on Conceptual Modeling (ER),
3 Springer, Brussels, LNCS, vol 6998, pp 104–117
- 4 15. Chopra AK, Artikis A, Bentahar J, Colombetti M, Dignum F, Fornara N, Jones AJI, Singh
5 MP, Yolum P (2013) Research directions in agent communication. *ACM Transactions on*
6 *Intelligent Systems and Technology (TIST)* In press
- 7 16. Desai N, Narendra NC, Singh MP (2008) Checking correctness of business contracts
8 via commitments. In: Proceedings of the 7th International Conference on Autonomous
9 Agents and MultiAgent Systems (AAMAS), IFAAMAS, Estoril, Portugal, pp 787–794
- 10 17. Falcone R, Castelfranchi C (2009) From dependence networks to trust networks. In:
11 Proceedings of the 11th AAMAS Workshop on Trust in Agent Societies (Trust), pp 13–26
- 12 18. Falcone R, Castelfranchi C (2010) Trust and transitivity: A complex deceptive relation-
13 ship. In: Proceedings of the 12th AAMAS Workshop on Trust in Agent Societies (Trust),
14 pp 43–54
- 15 19. Fornara N, Colombetti M (2002) Operational specification of a commitment-based agent
16 communication language. In: Proceedings of the 1st International Joint Conference on
17 Autonomous Agents and Multiagent Systems (AAMAS), ACM Press, Melbourne, pp
18 535–542
- 19 20. Fornara N, Colombetti M (2003) Defining interaction protocols using a commitment-
20 based agent communication language. In: Proceedings of the 2nd International Joint
21 Conference on Autonomous Agents and Multiagent Systems (AAMAS), ACM Press,
22 Melbourne, pp 520–527
- 23 21. Fornara N, Colombetti M (2009) Specifying and enforcing norms in artificial institutions.
24 In: *Declarative Agent Languages and Technologies VI*, Revised Selected and Invited
25 Papers, Springer, Berlin, LNCS, vol 5397, pp 1–17
- 26 22. Gasser L (1991) Social conceptions of knowledge and action: DAI foundations and open
27 systems semantics. *Artificial Intelligence* 47(1–3):107–138
- 28 23. Grosz B, Kraus S (1993) Collaborative plans for group activities. In: Proceedings of the
29 Twelfth International Joint Conference on Artificial Intelligence, pp 367–373
- 30 24. Habermas J (1984) *The Theory of Communicative Action*, volumes 1 and 2. Polity Press,
31 Cambridge, UK
- 32 25. Halpern JY, Moses YO (1990) Knowledge and common knowledge in a distributed
33 environment. *Journal of the Association for Computing Machinery* 37:549–587
- 34 26. Herrestad H, Krogh C (1995) Obligations directed from bearers to counterparties. In:
35 Proceedings of the 5th International Conference on Artificial Intelligence and Law, pp
36 210–218
- 37 27. Johnson M, Bradshaw JM, Feltovich PJ, Jonker CM, van Riemsdijk MB, Sierhuis M (2010)
38 The fundamental principle of coactive design: Interdependence must shape autonomy.
39 In: Proceedings of the AAMAS Workshop on Coordination, Organization, Institutions
40 and Norms (COIN), Springer, Toronto, LNCS, vol 6541, pp 172–191
- 41 28. Jones AJI, Parent X (2007) A convention-based approach to agent communication lan-
42 guages. *Group Decision and Negotiation* 16(2):101–141
- 43 29. Jones AJI, Sergot MJ (1993) On the characterisation of law and computer systems: the
44 normative systems perspective. In *Deontic Logic in Computer Science: Normative Sys-
45 tem Specification*. J. Wiley and Sons, 275–307
- 46 30. Levesque HJ, Cohen PR, Nunes JT (1990) On acting together. In: Proceedings of the
47 National Conference on Artificial Intelligence, pp 94–99
- 48 31. Marengo E, Baldoni M, Chopra AK, Baroglio C, Patti V, Singh MP (2011) Commitments
49 with regulations: Reasoning about safety and control in REGULA. In: Proceedings of
50 the 10th International Conference on Autonomous Agents and MultiAgent Systems
51 (AAMAS), IFAAMAS, Taipei, pp 467–474
- 52 32. McBurney P, Parsons S (2003) Dialogue game protocols. In: Huget MP (ed) *Communi-
53 cation in Multiagent Systems: Agent Communication Languages and Conversation
54 Policies*, LNAI, vol 2650, Springer, Berlin, pp 269–283

- 1 33. Minsky NH, Ungureanu V (2000) Law-governed interaction: A coordination and control
2 mechanism for heterogeneous distributed systems. *ACM Transactions on Software*
3 *Engineering and Methodology (TOSEM)* 9(3):273–305
- 4 34. Quine WvO (1960) *Word and Object*. MIT Press, Cambridge, MA
- 5 35. Sacerdoti E (1977) *The Structure of Plans and Behavior*. Elsevier North-Holland, New
6 York
- 7 36. Sichman JS, Conte R, Demazeau Y, Castelfranchi C (1994) A social reasoning mechanism
8 based on dependence networks. In: *Proceedings of the 11th European Conference on*
9 *Artificial Intelligence*, pp 188–192
- 10 37. Simon HA (1997) *Administrative Behavior: A Study of Decision-Making Processes in*
11 *Administrative Organizations*, 4th edn. Free Press, New York
- 12 38. Singh MP (1991) Social and psychological commitments in multiagent systems. In: *AAAI*
13 *Fall Symposium on Knowledge and Action at Social and Organizational Levels*, pp 104–
14 106
- 15 39. Singh MP (1998) Agent communication languages: Rethinking the principles. *IEEE Com-*
16 *puter* 31(12):40–47
- 17 40. Singh MP (1999) An ontology for commitments in multiagent systems: Toward a unifi-
18 cation of normative concepts. *Artificial Intelligence and Law* 7(1):97–113
- 19 41. Singh MP (2000) A social semantics for agent communication languages. In: *Proceedings*
20 *of the 1999 IJCAI Workshop on Agent Communication Languages*, Springer, Berlin,
21 *Lecture Notes in Artificial Intelligence*, vol 1916, pp 31–45
- 22 42. Singh MP (2007) Formalizing communication protocols for multiagent systems. In: *Pro-*
23 *ceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI)*,
24 *IJCAI, Hyderabad*, pp 1519–1524
- 25 43. Singh MP (2008) Semantical considerations on dialectical and practical commitments.
26 In: *Proceedings of the 23rd Conference on Artificial Intelligence (AAAI)*, AAAI Press,
27 Chicago, pp 176–181
- 28 44. Singh MP (2011) Trust as dependence: A logical approach. In: *Proceedings of the 10th*
29 *International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*,
30 *IFAAMAS, Taipei*, pp 863–870
- 31 45. Singh MP (2014) Norms as a basis for governing sociotechnical systems. *ACM Transac-*
32 *tions on Intelligent Systems and Technology (TIST)*, In press
- 33 46. Singh MP, Chopra AK, Desai N (2009) Commitment-based service-oriented architecture.
34 *IEEE Computer* 42(11):72–79
- 35 47. Telang PR, Singh MP (2011) Specifying and verifying cross-organizational business
36 models: An agent-oriented approach. *IEEE Transactions on Services Computing* 4, in
37 press
- 38 48. Udupi YB, Singh MP (2006) Contract enactment in virtual organizations: A commitment-
39 based approach. In: *Proceedings of the 21st National Conference on Artificial Intelligence*
40 *(AAAI)*, AAAI Press, Boston, pp 722–727
- 41 49. Udupi YB, Singh MP (2006) Multiagent policy architecture for virtual business organiza-
42 tions. In: *Proceedings of the 3rd IEEE International Conference on Services Computing*
43 *(SCC)*, IEEE Computer Society, Chicago, pp 44–51
- 44 50. Udupi YB, Singh MP (2007) Governance of cross-organizational service agreements:
45 A policy-based approach. In: *Proceedings of the 4th IEEE International Conference on*
46 *Services Computing (SCC)*, IEEE Computer Society, Salt Lake City, pp 36–43
- 47 51. Winikoff M, Liu W, Harland J (2005) Enhancing commitment machines. In: *Proceedings*
48 *of the 2nd International Workshop on Declarative Agent Languages and Technologies*
49 *(DALT)*, Springer, Berlin, LNAI, vol 3476, pp 198–220
- 50 52. Xing J, Singh MP (2003) Engineering commitment-based multiagent systems: A tem-
51 poral logic approach. In: *Proceedings of the 2nd International Joint Conference on*
52 *Autonomous Agents and MultiAgent Systems (AAMAS)*, ACM Press, Melbourne, pp
53 891–898
- 54 53. Yolum P, Singh MP (2002) Commitment machines. In: *Proceedings of the 8th Inter-*
55 *national Workshop on Agent Theories, Architectures, and Languages (ATAL 2001)*,
56 Springer, Seattle, LNAI, vol 2333, pp 235–247