# Computational-level Analysis of Constraint Compliance for General Intelligence

Robert E. Wray[0000−0002−5311−8593], Steven J. Jones[0000−0003−1942−6354], and John E. Laird[0000−0001−7446−3241]

The Center for Integrated Cognition
IQM Research Institute, Ann Arbor, MI 48105 USA
URL:integratedcognition.ai
{robert.wray,steven.jones,john.laird}@cic.iqmri.org

**Abstract.** Human behavior is conditioned by codes and norms that constrain action. Rules, "manners," laws, and moral imperatives are examples of classes of constraints that govern human behavior. These systems of constraints are "messy:" individual constraints are often poorly defined, what constraints are relevant in a particular situation may be unknown or ambiguous, constraints interact and conflict with one another, and determining how to act within the bounds of the relevant constraints may be a significant challenge, especially when rapid decisions are needed. General, artificially-intelligent agents must be able to navigate the messiness of systems of real-world constraints in order to behave predictability and reliably. In this paper, we characterize sources of complexity in constraint processing for general agents and describe a computational-level analysis for such *constraint compliance*. We identify key algorithmic requirements based on the computational-level analysis and outline a limited, exploratory implementation of a general approach to constraint compliance.

**Keywords:** Constraint compliance · Cognitive architecture

## 1 Introduction

Rules, social norms (e.g., "manners"), laws, and moral imperatives are examples of various classes of *constraints* that govern human behavior. Systems of constraints are "messy:" individual constraints are often poorly defined; the constraints relevant in a particular situation may be unknown or ambiguous; constraints interact and conflict with one another; and determining how to act rapidly within the bounds of relevant constraints may itself be a significant challenge. Yet humans routinely and robustly overcome the messiness of conforming to many simultaneous and often ill-defined constraints.

Notably, humans can also rapidly adapt their task performance to new constraints. A driver who has always driven on the left can, with just a little deliberation and practice, shift to driving on the right side of the road. A traveler has the ability to recognize and to adapt to overt local customs related to greetings,

meals, etc. Humans can quickly and robustly adapt to novel constraints, even when those novel constraints interact with familiar constraints and tasks.

Today's AI systems, in contrast, generally elide or ignore the messiness of complying with real-world constraints. They often encode a designer's interpretation of constraints (e.g., by knowledge engineering or learning from a human-defined policy) and are designed for limited, pre-specified operating contexts [13]. These systems conform to engineered constraints unfailingly but inflexibly. The encoding of constraints (along with designer assumptions) is tightly integrated with task specifications, making it difficult for the systems to adapt to new operating environments. For example, compare the relative immediacy of human adaptation to driving on their "opposite" side of the road for the first time vs. an autonomous driving system as trained today or the present limitations of large language models to conform to ethical guidance when producing responses [20].

These approaches can be acceptable for narrow AI but, as human intelligence suggests, a general artificial intelligence requires an ability to reason about its constraints (and conflicts), resolve ambiguity, determine how it should proceed given awareness of constraints, and be rapidly adaptive to new constraints. We introduce a broader approach to constraints, *constraint compliance*, intended to provide an agent with the capacity to comply with real-world constraints.

We consider the computational requirements for this more comprehensive approach to compliance to systems of constraints, emphasizing general intelligence. That is, we seek to identify a computational approach that is constraint-compliant, domain general (not specific to an application or a task domain), and robust to the complexities that "real world constraints" introduce. We outline sources of "messiness" relevant to constraint processing and present a computational analysis of an overall constraint-compliance process, enumerating five distinct types of processing steps the agent must make. We then outline an initial algorithmic-level exploration of a constraint-compliance process. Finally, based on the analysis and exploratory implementation, we identify four algorithmic challenges that require additional analysis and research in order to realize comprehensive constraint compliance.

## 2   Sources of Complexity in Constraint Processing

Here, we enumerate specific sources of complexity and challenge for comprehensive constraint compliance. This "messiness" derives from many sources spanning the environment, the agent's task(s), its internal capabilities and assumptions, and the specification of constraints themselves.

We illustrate using examples from Sudoku-puzzle solving and automobile driving. Sudoku is a canonical constraint satisfaction problem (CSP) [12] and offers an effective contrast between classical constraint satisfaction [4] and the more comprehensive account of constraint processing we examine here.[1]

---

[1] More recent approaches to constraints extend the coverage of classical approaches but do not span all the forms of messiness we consider [18].

Automobile driving offers a specific, familiar domain in which the real-world challenges of comprehensively complying to constraints arise; constraints abound in driving. This choice of domain is illustrative only: our goal is to develop a general approach to constraint compliance, not one specific to a single domain.

## 2.1  Partial Observability

In Sudoku, the puzzle state is fully available. The rules of the game (the constraints) can be readily applied after each move. Agents in real-world environments cannot generally sense everything and their actions often have uncertain outcomes. While partial observability and uncertainty have broad implications for agent reasoning [17], they impose specific demands for constraint compliance.

As one example, student drivers in the US are taught to "maintain a 3-second distance when following on dry payment." Unlike a speed limit, where a speedometer provides an immediate gauge of one's speed, fully complying with this constraint requires that the driver visually attend to and continually assess the distance and current speed of their vehicle vs. the one in front and adjust speed to maintain the minimum distance.[2] Because not all constraint-relevant parameters are directly accessible to the agent, the agent must take action to determine the compliance of its behavior with that constraint.

## 2.2  Dynamic, Fail-hard Environments

Dynamic environments compound the sources of messiness. Generally, dynamics amplifies the need for satisficing algorithmic solutions [19,8]. Algorithms must (minimally) be responsive to the dynamics of the environment. A driver cut off in traffic by another car cannot pause to reason about all the potential instantiations and implications of its constraints in this unexpected situation, it must continue to drive and manage its constraint compliance over time. The specification of constraints themselves can also change due to environment dynamics (e.g., new traffic laws). Finally, interactions between constraints (below) may only become evident as a dynamic situation unfolds.

## 2.3  Abstract and Poorly-defined Constraints

Real-world constraints are often ambiguous, abstract, and/or incomplete in their definition, giving rise to the challenge of interpreting and operationalizing such constraints [21]. In puzzles like Sudoku, however, constraint definition is unambiguous. Terms (cell, column, row) have immediate and direct correspondence to the representation of the puzzle. The constraints (or rules) defining the puzzle are also unambiguous.

In contrast, many constraints in driving are abstract ("drive defensively") or ambiguous ("do not follow too closely"). Terms used in constraints require a mapping onto one's internal representation that is not always consistent from person

---

[2] Some newer cars offer an indicator for travelling too closely. Thus, with a different embodiment, this constraint no longer requires active measurement.

to person. "Use caution near pedestrians" depends on how one understands and applies both "caution" and "near" and perhaps also "pedestrian."

It may seem possible to overcome this source of messiness by directly encoding the "meaning" of constraints into an agent. However, resilience and robustness in open-ended environments requires disintermediation of the encoding and interpretation of constraints. Attempting to specify in advance how the agent should interpret constraints in every situation is likely to fail when the agent (inevitably) encounters a situation not anticipated by a system designer.

### 2.4   Implicit Context Specification

The definitions of real-world constraints often imply additional parameters or conditions rather than explicitly defining them. Most importantly, constraint specifications typically omit the context(s) in which they should apply. By "context," we mean a set of situations that share common, salient features. The "automobile driving" context includes cars, roads, traffic laws, traffic signals, etc. Similar but different contexts can have constraints that prescribe very different behaviors. For instance, "do not pass on the right" is a constraint relevant in countries where vehicles are driven on the right side of the road, but is not apt (most of the time) for countries where vehicles are driven on the left.

For Sudoku, there is an implicit but single context. Thus implicit specification poses no problem to the classical approach to constraints.

### 2.5   Interactions & Conflicts among Constraints, Tasks, & Contexts

Interactions and conflicts among constraints and between task(s) and constraints can arise frequently. An accident or road construction causes re-routing of traffic into normally oncoming traffic lanes. A text message notification draws attention when attending to the road is required (sometimes by law). To what extent should one obey traffic laws when transporting someone in dire medical distress? The specific instantiation of constraints grounded within a given situation will indicate competing and sometimes conflicting choices for the agent.

The design of Sudoku ensures that constraints are collectively coherent (simplified by the single context). Generally, classical approaches to satisfying constraints only provide solutions when sets of constraints are coherent, obviating conflicts. More recent approaches support "soft constraints" [15] which support prioritization of constraints when conflicts arise; the overall set of constraints remains coherent when prioritization is taken into account.

In real-world situations, conflicts cannot always be resolved via *a priori* prioritization; an agent must sometimes knowingly violate a constraint. If a car is cut off in heavy traffic, it is probably more important to maintain speed and slowly build distance between the car ahead than to sharply brake in order to regain compliance with the following-distance constraint. From the point of view of constraint compliance, the agent is often likely to be in situations, imposed by dynamics directly but also sometimes at its choosing given the dynamics, to violate some constraints and to repair violations as the evolving situation allows.
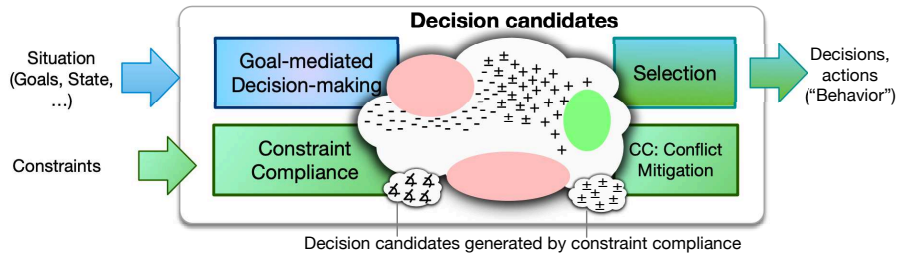
Decision candidates generated by constraint compliance

**Fig. 1.** At the computational level, the purpose of constraint compliance is to ensure that decision making takes constraints relevant to the agent's situation into account.

## 3   Computational-level Analysis

We now present a computational-level account of the information processing tasks necessary for constraint compliance given the many sources of "messiness" above. A computational-level analysis emphasizes *what* steps are required to achieve constraint compliance and identifies requirements for *how* the capability may be realized at the processing and representation ("algorithmic") level [14].

### 3.1   Functional Role

The functional role of constraint compliance is to modulate agent decisions (and thus behavior) so that constraints relevant to the current situation inform agent choices. In Figure 1, the agent's goal-focused decision process (blue) generates candidate choices and selects among them. The primary input to this decision process is the current situation (including environment state, external goals, history, etc.) and the output is a decision. A decision could be a commitment to a long-term course of action (e.g., a plan), an intermediate subgoal, or an immediate action. Over time, the sequence of decisions produces behavior (e.g., "driving"). We illustrate constraint compliance (green) parallel to the goal-mediated decision process of the agent and external constraints as a distinct input. This separation is for illustration only; at the algorithmic level, solutions may integrate constraint-compliance with goal-mediated decision processes.

As suggested by the figure, the agent commits to its decisions from a (potentially very large) space of candidate choices. At the computational level, we do not assume that the agent has an explicit representation of this space; in the figure, the cloud represents a conceptual space from which a specific decision might be drawn. For example, the decision process might choose actions based on a learned policy, where the space is implicit in mappings from states to actions.

Functionally, constraint compliance augments candidate choices produced by the goal-mediated decision process by indicating the acceptability/desirability of the candidates with respect to relevant constraints. The figure shows parts of the candidate space that are required (green), prohibited (red), desirable (+), and undesirable (-) choices. Because constraints can conflict (§2.5), some candidates
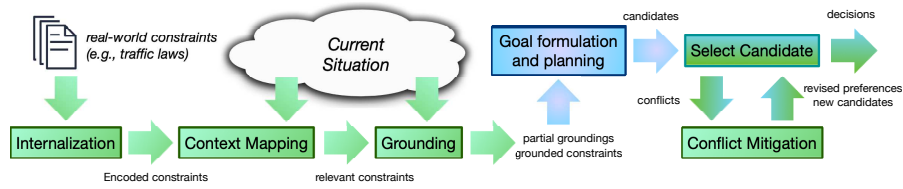
**Fig. 2.** Simple process model for constraint compliance.

are labeled as both desired and undesired ($\pm$); however, conflicts can occur in any combination. The selection process (green/blue) now evaluates the candidates and the desirability of those candidates.

Constraint compliance can also add new candidates. The grounding process can suggest candidates to take new actions (e.g., measurements to evaluate individual constraints; §2.1). In order to mitigate conflicts in constraints, the selection process may produce new candidates as well. Thus, in contrast to classical constraint satisfaction (where the application of constraints reduces choices), constraint compliance can produce additional choices. It also enables the agent to choose courses of action that are not necessarily consistent with all constraints.

### 3.2   Processing Steps for Constraint Compliance

What computational tasks are performed by the constraint-compliance process? Figure 2 illustrates a high-level process. The specific sequence of steps illustrates both a simple process model and how we are exploring constraint compliance at the algorithmic level and integrating it with decision making (see §4).

The agent's internal representations of constraints derive from real-world constraints defined externally (e.g., a law). *Internalization* results in encoding of constraints in agent memory. Next, *Context Mapping* compares encoded constraints to the current situation, identifying what constraints are (potentially) relevant in a given situation. Context mapping results in a set of situation-relevant but abstract (not grounded) constraints.

*Grounding* then maps abstract constraints to specific objects in the environment. In our explorations to date, both complete and partially-grounded constraints are re-represented as goals in order to exploit an existing agent's planning capability (§4). Planning generates candidates for *Selection* which is now extended with an ability to assess the acceptability of decision candidates based on the constraints. When conflicts arise, selection is augmented with *Conflict Mitigation*, which may lead to the generation of of alternative courses of action. Below, we further describe these steps, focusing especially on how "messiness" motivates and/or introduces additional requirements for individual steps.

**Internalize constraints.** Real-world constraints (typically) are defined external to the agent. Thus, an initial step in constraint compliance is to interpret the external constraint; that is, to map the external representation of the constraint to concepts as represented within the agent.

Abstract and poorly-defined constraints (§2.3) introduce challenges to simple encoding. The agent may not possess internal representations that align with the conditions in the external constraint and thus algorithmic approaches to internalization will entail methods that allow an agent to assess mappings between external conditions and internal representations.

**Identify situational context(s).** Conforming to real-world constraints requires an agent to recognize which constraints are relevant to its situation. However, the applicable situation (or general characterization of situations: contexts) are often implicit in the specification of constraints (§2.4). An agent can often learn associations between contexts and constraints through experience (which can include instruction) but a core challenge is that constraint specifications themselves do not (usually) specify applicable contexts.

A second challenge results when the composition of contexts interact in ways that make previously learned mappings inapt or invalid (§2.5). Anticipating and evaluating all possible compositions of all possible contexts is not feasible. Thus, general intelligence requires the capacity to consider and to evaluate constraints in novel contexts as behavior is being generated.

Context recognition itself is a challenge [6]. For an algorithmic implementation of constraint compliance, all that is needed is that the agent recognize "this constraint is relevant in my current situation." However, a full solution to constraint compliance appears to require context recognition processes as well.

**Instantiate constraints in a situation (Grounding).** As an agent behaves in its environment, it must determine how constraints might apply in its current situation. Grounding is distinct from internalization and context identification; it requires that the agent shift from general consideration of a constraint to determining if/how it should be instantiated in the agent's current environment.

Grounding is often straightforward. However, partial observability (§2.1) and abstract constraints (§2.4) can require a search over potential instantiations of a constraint, rather than an immediate mapping. Thus, as constraints are expressed more abstractly and generally, the computational demand on the agent to determine *how* that constraint may apply in the current situation increases. When new information is needed to complete grounding (e.g., a measurement as in §2.1), new candidate choices should be generated (§3.1).

An agent's embodiment may lack an ability to directly observe features needed to instantiate a constraint. Nonetheless, the agent should still attempt to respect applicable constraints. Thus, grounding requires prospective instantiation with incomplete information.

**Integrate constraints in decision-making (Selection).** At a minimum, the agent's selection process must take into account both agent goals and constraints for constraint compliance. When the set of applicable constraints are fully grounded and present no conflicts, the selection process is straightforward.

Conflicts (below) and partial grounding complicate selection. The selection process must be sensitive to both taking action to find an instantiation for a partially grounded constraint and also the potential costs and risks associated

with that search. Defining algorithmic approaches to selection in the presence of partial grounding is a significant novel challenge.

**Identify and mitigate conflicts.** When there are conflicts in the acceptability and desirability of candidate choices, the agent must either 1) choose one of the options given the conflicting choices or 2) attempt to identify new choices that resolve or mitigate the conflicts. Specific strategies could include attention/inattention (ignoring some constraints), prioritization of constraints, and replanning. A primary algorithmic-level challenge is to resolve and mitigate conflicts rapidly, given bounded rationality in a dynamic environment (§2.2).

## 4   Exploratory Algorithmic-level Prototype

In parallel with the top-down computational-level analysis, we have begun bottom-up prototyping as well, focusing to date on algorithmic approaches to grounding, selection, and conflict mitigation. We use Soar [11] as the target implementation level. Soar both constrains and informs definition at the algorithmic level. We introduce further design constraint at the algorithmic level by building on an existing agent designed to interactively learn tasks [10,16]. The prototype is compatible with this agent's *a priori* capabilities for interpreting language, planning task actions, executing plans, and learning from human instruction.

**Grounding:** The prototype builds on language grounding already part of the agent, which can learn recognition structures for abstract goal specifications [10] and maintain consistent grounding across perceptual changes [16]. The primary focus is to explore how to support partial grounding of constraints. The agent can now indicate that some actions are desirable (in Selection; see below) because they lead to further information that could potentially complete the grounding. In this way, the agent is biased towards choosing candidates that lead to measurement actions, as suggested in Figure 1.

**Selection:** The original agent uses an explicit goal representation to determine what to do next (typically via search-based planning, although it can ask for help from an instructor as well). In our initial implementation, as shown in Figure 2, we integrated constraint-compliance with selection by having the agent represent grounded constraints as goals (e.g., a speed limit constraint would be represented as a goal for speed to be less than the limit). This approach leverages the agent's planning capability. Candidate evaluations (from grounding) are implemented as Soar preferences for selecting plans, which maps selection directly onto an implementation/architecture-level capability of Soar. In the absence of conflicts (below), planning provides a solution that satisfies the (grounded) constraints, measurement actions (from partial groundings), and task actions.

**Conflict Mitigation:** Consider two conflicting constraints relevant to driving in a medical emergency. The lawful speed limit and a general directive to preserve human life apply. These constraints can result in a conflict over the desired speed. Because plan choices are mapped onto Soar preferences, Soar responds to conflicting preferences with an impasse, a conflict detection system already part of Soar. Thus, we have also mapped the trigger for conflict mitiga-

tion onto an implementation-level process. Generally, resolving conflicts requires additional knowledge (e.g., in this case, some sense that preserving life is more important than respecting the speed limit) which can include various ways to include values in assessing choices [1,7].

# 5  Discussion and Implications

While limited and preliminary, the initial prototype highlights examples of representation and process (algorithmic-level choices) and how these choices may interact with the implementation level. We now consider implications for future work at the algorithmic level to realize general constraint compliance.

**Online, Incremental Learning:** For an AGI, the set of contexts and constraints is potentially huge, it is infeasible to prepare for every contingency, and dynamics often demands rapid response. Together, these conditions point toward algorithmic solutions that employ online, incremental learning. This implication mirrors human learning and is consistent with the transition from more deliberate and explicit (System 2) to more implicit and automatic (System 1) reasoning [9]. However, it contrasts with recent approaches that emphasize pre-training to ensure conformance to various operational and safety constraints [5].

**Senses of Familiarity, Novelty, and Surprise:** Familiarity, novelty, and surprise are important signals in human (and animal) regulation of behavior [2]. Realizations of familiarity, novelty and surprise may be useful for meta-cognitive regulation of constraint compliance in task performance. An open question is whether a sense of familiarity (and other signals) are best realized in the implementation level (e.g., extension to Soar) or at the algorithmic level.

**Anticipation based on Partial Information:** Near-term anticipation of future states is central to functional and neurological accounts of human intelligence [3]. Humans readily anticipate the potential impact of constraints on behavior and adapt behavior in advance of a potential constraint violation. Our exploration identified a need for anticipation in grounding. An agent needs strategies to decide which potential groundings to attend to, given many potential groundings (with many implications). Indicators of potential threats to constraint compliance would provide a coarse attention mechanism to bias grounding processes toward more important constraints.

**Domain Knowledge:** Choosing to prioritize some constraints over others requires general knowledge of the world. Having such knowledge may be as important to the results of constraint compliance as the algorithms that realize its functions. This dilemma points to one of the rationales for adopting an agent that can learn from instruction. Because research agents will often lack knowledge, our agent can actively seek input to gain missing knowledge about conflicts. While this does not resolve the dependence of constraint compliance on general knowledge, it does provide a means to explore algorithmic realizations in a way that makes the required domain knowledge explicit and transparent.

# References

1. Arkin, R.C., Ulam, P., Wagner, A.R.: Moral decision making in autonomous systems: Enforcement, moral emotions, dignity, trust, and deception. Proceedings of the IEEE **100**(3), 571–589 (2011)
2. Barto, A., Mirolli, M., Baldassarre, G.: Novelty or Surprise? Front. in Psy. **4** (2013)
3. Bubic, A., von Cramon, D.Y., Schubotz, R.I.: Prediction, Cognition and the Brain. Frontiers in Human Neuroscience **4**,  25 (Mar 2010)
4. Dechter, R.: Constraint processing. Morgan Kaufman (2003)
5. García, J., Fern, Fernández, O.: A Comprehensive Survey on Safe Reinforcement Learning. Journal of Machine Learning Research **16**(42), 1437–1480 (2015)
6. Gershman, S.J.: Context-dependent learning and causal structure. Psychonomic bulletin & review **24**, 557–565 (2017)
7. Giancola, M., Bringsjord, S., Govindarajulu, N.S., Varela, C.: Ethical reasoning for autonomous agents under uncertainty. In: International Conference on Robot Ethics and Standards (ICRES). pp. 1–16. Taipei (2020)
8. Gigerenzer, G.: Fast and frugal heuristics: Tools of bounded rationality. In: Handbook of judgment and decision making, pp. 62–88. Blackwell, Malden (2004)
9. Kahneman, D.: Thinking, fast and slow. Doubleday, New York (2011)
10. Kirk, J.R., Laird, J.E.: Learning Hierarchical Symbolic Representations to Support Interactive Task Learning and Knowledge Transfer. In: IJCAI 2019. IJCAI (2019)
11. Laird, J.E.: The Soar Cognitive Architecture. MIT Press, Cambridge, MA (2012)
12. Lynce, I., Ouaknine, J.: Sudoku as a SAT Problem. In: AI&M (2006)
13. Mani, G., Chen, F., et al.: Artificial Intelligence's Grand Challenges: Past, Present, and Future. AI Magazine **42**(1), 61–75 (Apr 2021), number: 1
14. Marr, D.: Vision. Freeman and Company, New York (1982)
15. Meseguer, P., Rossi, F., Schiex, T.: Soft constraints. In: Foundations of Artificial Intelligence, vol. 2, pp. 281–328. Elsevier (2006)
16. Mininger, A.: Expanding Task Diversity in Explanation-Based Interactive Task Learning. Ph.D. Thesis, University of Michigan, Ann Arbor (2021)
17. Pearl, J.: Reasoning Under Uncertainty. Annual Review of Computer Science **4**(1), 37–72 (1990), _eprint: https://doi.org/10.1146/annurev.cs.04.060190.000345
18. Rossi, F., Mattei, N.: Building ethically bounded AI. In: $33^{rd}$ AAAI Conf. (2019)
19. Simon, H.A.: Models of man; social and rational. Wiley, Oxford, England (1957)
20. Weidinger, L., Mellor, J., et. al: Ethical and social risks of harm from Language Models (Dec 2021), arXiv:2112.04359 [cs]
21. Wray, R.E., Laird, J.E.: Incorporating Abstract Behavioral Constraints in the Performance of Agent Tasks. In: ICAI. Springer, Las Vegas, NV (Jul 2021)