# Liability, Ethics, and Culture-Aware Behavior Specification using Rulebooks

Andrea Censi, Konstantin Slutsky, Tichakorn Wongpiromsarn,
Dmitry Yershov, Scott Pendleton, James Fu, Emilio Frazzoli

*Abstract*— The behavior of self-driving cars must be compatible with an enormous set of conflicting and ambiguous objectives, from law, from ethics, from the local culture, and so on. This paper describes a new way to conveniently define the desired behavior for autonomous agents, which we use on the self-driving cars developed at nuTonomy, an Aptiv company.

We define a "rulebook" as a pre-ordered set of "rules", each akin to a violation metric on the possible outcomes ("realizations"). The rules are ordered by priority. The semantics of a rulebook imposes a pre-order on the set of realizations. We study the compositional properties of the rulebooks, and we derive which operations we can allow on the rulebooks to preserve previously-introduced constraints.
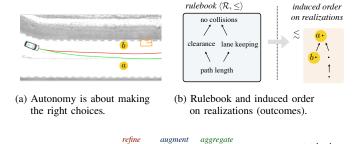
While we demonstrate the application of these techniques in the self-driving domain, the methods are domain-independent.

## I. INTRODUCTION

One of the challenges in developing self-driving cars is simply *defining what the car is supposed to do*. The behavior specification for a self-driving car comes from numerous sources, including not only the vaguely specified "rules of the road", but also implementation limitations (for example, the speed might be limited due to the available computation for perception), and numerous other soft constraints, such as the need of "appearing natural", or to be compatible with the local driving culture (any reader who never lived in Boston will be surprised to discover what is the "Massachusetts left"). As self-driving cars are potentially life-endangering, also moral and ethical factors play a role [1]–[3]. For a self-driving car, the "trolley problems" [4]–[6] are not idle philosophical speculations, but something to solve in a split second. As of now, there does not exist a formalism that allows to incorporate all these factors in one specification, which can be precise enough to be taken as regulation for what self-driving cars designers must implement.

Formal methods have been applied to specify and verify properties of complex systems. The main focus has been to provide a proof that the system satisfies a given specification, expressed in a formal language. In particular, specifications written in temporal logics have been studied extensively [7]–[12]. In self-driving cars, often times, not all the rules can be satisfied simultaneously. Although there are formalisms that allow specifying the degree of satisfaction of each rule, e.g., based on fuzzy logic or some measurable probability [13], [14], as of now, there does not exist a formalism that allows to incorporate different factors needed to be considered for self-driving cars with a precise hierarchy in one specification.

The authors are with nuTonomy, an Aptiv company (Boston, MA, Zurich, and Singapore). Please address correspondence to andrea@nutonomy.com.

(a) Autonomy is about making the right choices.

(b) Rulebook and induced order on realizations (outcomes).

(c) Rulebook manipulation operations refine the specification

Fig. 1: The rulebooks formalism allows to specify the desired behavior for an autonomous agent by using a pre-ordered set of rules that induce a pre-order on the allowed outcomes. The rulebooks can be refined by a series of manipulation operations.

In this paper we describe a formalism called "rulebooks", which we use to specify the desired behavior of the self-driving cars developed at nuTonomy[1]. While the formalism can be applied to any system, it is particularly well-suited to handle behavior specification for embodied agents in an environment where many, possibly conflicting rules need to be honored. We define a "rulebook" as a set of "rules" (Fig. 1b), each akin to a violation metric on the possible outcomes. The rules can be defined analytically, using formalisms such as LTL [15] or even deontic logic [16], or the violation functions can be learned from data, using inverse reinforcement learning [17], or any technique that allows to measure deviation from a model. In the driving domain, the rules can derive from traffic laws, from common sense, from ethical considerations, etc.

The rules in a rulebook are hierarchically ordered to describe their relative priority, but, following the maxim "*good specifications specify little*", the semantics of a rulebook imposes a *pre-order* on the set of outcomes, which means that the implementations are left with considerable freedom of action. The rulebooks formalism is "user-oriented": we define a set of intuitive operations that can be used to iteratively

---

[1]Please note that the functionality described is not necessarily representative of current and future products by nuTonomy, Aptiv and their partners. The scenarios discussed are simplified for the purposes of exposition. The specification examples discussed are illustrative of the philosophy but not the precise specification we use. The methodology described does not represent the full development process; in particular we gloss over the extensive verification and validation processes that are needed for safety-critical rules.

refine the behavior specification. For example, one might define an "international rulebook" for rules that are valid everywhere, and then have region-specific rulebooks for local rules, such as which side the car should drive on.

While the rulebooks offer formidable generality in describing behavior, at the same time, when coupled with graph-based motion planning, the rulebooks allow a *systematic, simple, and scalable solution to planning* for a self-driving car:

1) Liability-, ethics-, culture-derived constraints are formulated as rules (preferences over trajectories), either manually or in a data-driven fashion, together with the rules of the road and the usual geometric constraints for motion planning.
2) Priorities between conflicting rules are established as a rulebook (ideally, by nation-wide regulations based on public discourse);
3) Developers customize the behavior by resolving ambiguities in the rulebook until a total order is obtained;
4) Graph-based motion planning, in particular variations of minimum-violation motion planning [18]–[22], allow to generate the trajectories that maximally respect the rules in the rulebooks.

In a nutshell, the above is how the nuTonomy cars work. The topic of *efficiently* planning with rulebooks is beyond the goals of this paper; here, we focus on the use of the rulebooks as a specification, treating the planning process as a black box.

## II. RULEBOOKS DEFINITION

*1) Realizations:* Our goal is to define the desired agent *behavior*. Here, we use the word "behavior" in the sense of Willems [23] (and not in the sense of "behavior-based robotics" [24], [25]), to mean that what we want to prescribe is what we can measure objectively, that is, the externally observable actions and outcomes in the world, rather than the internal states of the agent or any implementation details. Therefore, we define preference relations on a set of possible outcomes, which we call the set of *realizations* $\Xi$. For a self-driving car, a realization $x \in \Xi$ is a world trajectory, which includes the trajectory of all agents in the environment.

We use *no* concept of infeasibility. Sometimes the possible outcomes are all catastrophically bad; yet, an embodied agent must keep calm and choose the least catastrophic option.

*2) Rules:* Our "atom" of behavioral specification is the "rule". In our approach, a rule is simply a scoring function, or "violation metric", on the realizations.

**Definition 1** (Rule). Given a set of realizations $\Xi$, a *rule* on $\Xi$ is a function $r : \Xi \to \mathbb{R}_+$.

The function $r$ measures the degree of violation of its argument. If $r(x) < r(y)$, then the realization $y$ violates the rule $r$ to a greater extent than does $x$. In particular, $r(x) = 0$ indicates that a realization $x$ is fully compliant with the rule.

Any scalar function will do. The definition of the violation metric might be analytical, "from first principles", or be the result of a learning process.

In general, the rulebooks philosophy is to pay particular attention about specifying what we ought to do *when the rule has to be violated*, as described in the following examples.

**Example 2** (Speed limit). A naïve rule that is meant to capture a speed limit of 45 km/h could be defined as:

$$r(x) = \begin{cases} 0, & \text{if the car's speed is always below } 45\,\text{km/h}, \\ 1, & \text{otherwise.} \end{cases}$$

However, this discrete penalty function is not very useful in practice. The rulebooks philosophy is to assume that rules might need to be violated for a greater cause. In this case, it is advisable to define a penalty such as:

$$r(x) = \text{interval for which the car was above 45 km/h.}$$

The effect of this will be that the car will try to stay below the speed limit, but if it cannot, it will minimize the time spent violating the limit. Alternatively, one can penalize also the magnitude of the speed violation:

$$r'(x) = r(x) \times (v_{\max} - 45 \text{ km/h}).$$

**Example 3** (Minimizing harm). It is easy enough to write a constraint describing the fact that we do not want any collision; but, assuming that a collision with a human is unavoidable given the circumstances, what should the car do? In this case, it would be advisable to define the violation function as:

$$r(x) = \text{kinetic energy transferred to human bodies,}$$

so that the car will try to avoid collisions, but, if a collision is inevitable, it will try to reduce the speed as much as possible.

*3) Rulebooks:* A rulebook $\mathcal{R}$ is a *pre-ordered* set of rules. We will use $\mathcal{R}$ both for the rulebook and for its underlying set of rules.

**Definition 4** (Rulebook). A *rulebook* is a tuple $\langle \mathcal{R}, \leq \rangle$, where $\mathcal{R}$ is a finite set of rules and $\leq$ is a preorder on $\mathcal{R}$.
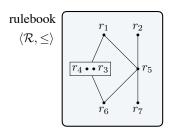


Fig. 2: Graphical representation of a rulebook. Rules are ordered vertically with the most important rules being at the top.

Being a preorder, any rulebook may be represented as a directed graph, in which each node is a rule, and an edge between two rules $r_1 \to r_2$ means that $r_1 \leq r_2$, i.e., the rule $r_2$ has higher rank. Fig. 2 gives an example of a rulebook with 7 rules. In this example, rules $r_1$ and $r_2$ are incomparable, but both are greater than $r_5$. Rules $r_3$ and $r_4$ are of the same rank, meaning $r_3 \leq r_4$ *and* $r_4 \leq r_3$, and both are smaller than $r_1$, greater than $r_6$ and incomparable to $r_5$, $r_2$, or $r_7$.

Just like it might be convenient to learn some of the non-safety-critical rules from data, it is possible to learn some of the priorities from data as well. (See [26], [27] for a similar concept in a different context.)

*4) Induced pre-order on realizations:* We now formally define the semantics of a rulebook as specifying a pre-order on realizations. Because a rulebook is defined as a *pre-ordered* set of rules, not all the relative priorities among different rules are specified. We will see that this means that a rulebook can be used as a very flexible *partial* specification.

Given a rulebook $\langle \mathcal{R}, \leq \rangle$, our intention is to preorder all realizations such that $x \lesssim y$ can be interpreted as $x$ being "at least as good as" $y$, i.e., the degree of violation of the rules by $x$ is at most as much as that of $y$.

**Definition 5** (Pre-order $\lesssim$ and strict version $<$ ). Given a rulebook $\langle \mathcal{R}, \leq \rangle$ and two realizations $x, y \in \Xi$, we say that $x \lesssim y$ if for any rule $r \in \mathcal{R}$ satisfying $r(y) < r(x)$ there exists a rule $r' > r$ such that $r'(x) < r'(y)$. We denote by $<$ the strict version of $\lesssim$.

**Lemma 6.** *Let $\langle \mathcal{R}, \leq \rangle$ be a rulebook, let $x, y, z \in \Xi$ be realizations such that $x \lesssim y$, $y \lesssim z$, and let $r \in \mathcal{R}$ be a rule. If either $r(x) \neq r(y)$ or $r(y) \neq r(z)$, then there exists a rule $r' \geq r$ such that $r'(x) < r'(z)$.*

*Proof.* We give a proof for $r(x) \neq r(y)$; the case $r(y) \neq r(z)$ is analogous. If $r(x) > r(y)$, then $x \lesssim y$ guarantees existence of $r_0 > r$ such that $r_0(x) < r_0(y)$. If $r(x) < r(y)$ to begin with, then we may set $r_0 = r$, and in either case we get $r_0 \geq r$ such that $r_0(x) < r_0(y)$. We are done if $r_0(y) \leq r_0(z)$, as one can take $r' = r_0$. Otherwise, $y \lesssim z$ implies existence of $r_1 > r_0$ such that $r_1(y) < r_1(z)$. Again, we are done if $r_1(x) \leq r_1(y)$, and if not, there has to be some rule $r_2 > r_1$ such that $r_2(x) < r_2(y)$. Continuing in the same fashion, one builds an increasing chain $r \leq r_0 < r_1 < r_2 < \cdots$. Since $\mathcal{R}$ is assumed to be finite, the chain has to stop, which is possible only if $r_n(x) < r_n(z)$ for some $n$.

**Proposition 7.** *Let $\langle \mathcal{R}, \leq \rangle$ be a rulebook, and let $x, y, z \in \Xi$ be realizations.*
*1) The relation $\lesssim$ on realizations is a preorder.*
*2) Two realizations $x$ and $y$ are equivalent if and only if $r(x) = r(y)$ for all rules $r \in \mathcal{R}$.*

*Proof.* 1) It is clear that $\lesssim$ is reflexive, so we only need to check transitivity. Suppose $x \lesssim y$, $y \lesssim z$, and let $r \in \mathcal{R}$ be such that $r(x) > r(z)$. Clearly either $r(y) \neq r(x)$ or $r(z) \neq r(y)$, so Lemma 6 applies, producing $r' > r$ such that $r'(x) < r'(z)$, hence $x \lesssim z$ as claimed.

2) Suppose towards a contradiction there are some realizations satisfying $x \lesssim y$ and $y \lesssim x$, yet $r(x) \neq r(y)$ for some $r \in \mathcal{R}$. Without loss of generality, let us assume that $r(x) < r(y)$. Since we have $x \lesssim y \lesssim x$, Lemma 6 produces some $r' \geq r$ such that $r'(x) < r'(x)$, which is absurd. The other direction ($\forall r, r(x) = r(y) \implies x \sim y$) is obvious.

*Remark 8.* In the special case in which the rulebook is a linear order, the induced order on realizations is the *lexicographic order* used in the literature in minimum-violation planning.

## III. EXAMPLES IN THE DRIVING DOMAIN

In this section, we give a few examples of the types of rules that are useful in the driving domain. Rather than describing the full complexity of our production rules, which address subtle nuances of behavior and idiosyncrasies and corner cases of traffic laws, we prefer to give a few synthetic examples of rulebooks and rulebooks refinement.

**Example 9** (Safety vs. infractions)**.** Consider the scenario in Fig. 3. A vehicle is faced with an obstacle in front, and is given a choice between two trajectories $a$ and $b$. Suppose the initial speed of the vehicle is sufficiently high, and there is no time to stop, so collision is unavoidable if $a$ is chosen. Trajectory $b$, however, is collision free, but it violates a different rule, since it intersects a double solid line.

The rulebooks take on this situation is the following. A rule "not to collide with other objects" will have a higher priority than the rule of not crossing the double line (Fig. 3b). With this rulebook, the trajectory $b$ will be chosen to avoid the collision.
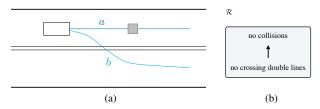


Fig. 3: The rulebook allows the agent to cross the double white line to avoid a collision. (This assumes that there are no other agents outside the frame that might trigger the "no-collision" rule.)
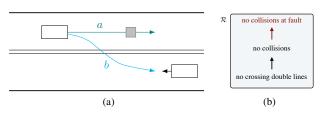


Fig. 4: The rulebook instructs the agent to collide with the object on its lane, rather than provoking an accident, for which it would be at fault.

**Example 10** (Liability-aware specification)**.** Let's change the situation slightly by assuming that trajectory $b$ is also in collision, but with a different agent — an oncoming vehicle on the opposite lane. Under these assumptions, we may be interested in choosing the outcome where the ego vehicle is *not* at fault for the collision.

This behavior specification can be achieved by the rulebook of Fig. 4b, having two collision rules, where one evaluates the degree of collision, where the ego vehicle is at fault, and the other evaluates collisions caused by third-party, which is below the former in the rulebooks hierarchy. This will force the ego vehicle to prefer trajectory $a$ over $b$.

This example fully captures the concept of the "responsibility-sensitive safety" model described in [28].

**Example 11** (Partial priorities specification). Consider the scenario depicted in Fig. 5a, where the vehicle encounters an obstacle along its route. For simplicity, we focus on four discrete representative trajectories, called $a, b, c, d$. A minimal rulebook that allows to deal with this situation would contain at least four rules, detailed below. For simplicity, we write the violation metrics as binary variables having value 0 or 1 on the test trajectories, while in practice these would be continuous functions.

1) Rule $\beta$ - **Blockage**, attaining value 1 if the trajectory is blocked by an obstacle, and 0 otherwise:

$$\beta(x) = \begin{cases} 0, & \text{for } x = b, c, d; \\ 1, & \text{for } x = a. \end{cases}$$

2) Rule $\lambda$ - **Lane Keeping**, 1 iff the trajectory intersects the lane boundary:

$$\lambda(x) = \begin{cases} 0, & \text{for } x = a, b; \\ 1, & \text{for } x = c, d. \end{cases}$$

3) Rule $\kappa$ - **Obstacle clearance**, 1 iff the trajectory comes closer to an obstacle than some threshold $C_0$:

$$\kappa(x) = \begin{cases} 0, & \text{for } x = c, d; \\ 1, & \text{for } x = a, b. \end{cases}$$

*Remark* 12 (Learning while preserving safety). While parameters such as the minimum clearance from an obstacle $C_0$ can be specified manually, in practice, given an adequate data analytics infrastructure, they are great candidates to be *learned* from the data. This allows the car to adapt the behavior to the local driving culture. By still having the safety-preserving rules at the top of the hierarchy, the rulebooks allow the system to be adaptive without ever compromising safety, not even with adversarial data. (See [29] for a similar principle in a different context.)
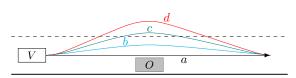
4) Rule $\alpha$ - **Path length**, whose value is the length of the trajectory:

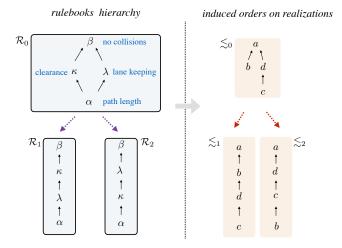$$\alpha(a) < \alpha(b) < \alpha(c) < \alpha(d).$$

Out of these rules, we can make different rulebooks by choosing different priorities. For example, defining the rulebook $\mathcal{R}$ with ordering $\alpha < \kappa < \beta$ and $\alpha < \lambda < \beta$, depicted in Fig. 5b, the following order on trajectories is imposed: $b < a$ and $c < d < a$. Note that $b$ is not comparable with either $d$ or $c$. This is an important feature of a *partial specification:* we leave freedom to the implementation to choose the details of the behavior that we do not care about.

## IV. ITERATIVE SPECIFICATION REFINEMENT WITH RULEBOOKS MANIPULATION

We formalize this process of *iterative specification refinement* (Fig. 1c), by which a user can add rules and priority relations until the behavior of the system is fully specified to one's desire.



(a) Trajectories available to a vehicle before an avoidance maneuver.



(b) Rulebook hierarchy and induced hierarchy on realizations order.

Fig. 5: Example involving an avoidance maneuver.

**Example 13.** Regulations in different states and countries often share a great deal of similarity. It would be ineffective to start the construction of rulebooks from scratch in each case; rather, we wish to be able to define a "base" rulebook that can then be particularized for a specific legislation by adding rules or priority relations.

*1) Operations that refine rulebooks:* We will consider three operations (Fig. 6):

1) **Priority refinement** (Def. 14): this operation corresponds to adding another edge to the graph, thus clarifying the priority relations between two rules.
2) **Rule aggregation** (Def. 16): this operation allows to "collapse" two or more equi-ranked rules into one.
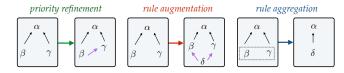3) **Rule augmentation** (Def. 17): this operation consists in adding another rule at the lowest level of priority.



Fig. 6: Three operations for manipulation of rulebooks.

*2) Priority refinement:* The operation of refinement adds priority constraints to the rulebook.

**Definition 14.** An allowed *priority refinement* operation of a rulebook $\langle \mathcal{R}_1, \leq_1 \rangle$ is a rulebook $\langle \mathcal{R}_1, \leq_2 \rangle$, where the order $\leq_2$ is a refinement of $\leq_1$.

**Example 15.** Continuing the example in Fig. 5, we can create two refinements of the rulebook by adding priority constraints

that resolve the incomparability of rules $\kappa$ and $\lambda$ one way or the other. For example, choosing the totally ordered rulebook $\alpha \to \kappa \to \lambda \to \beta$, the order on trajectories is $b < c < d < a$, while for the rulebook $\alpha \to \lambda \to \kappa \to \beta$, the order is $c < d < b < a$.

*3) Rule aggregation:* Suppose that a rulebook includes two rules that are in the same equivalence class. The minimal example is a rulebook $\langle \mathcal{R}, \leq \rangle$ that has two rules $r_1, r_2$ such that $r_1 \leq r_2$ and $r_2 \leq r_1$. The induced order $\lesssim$ on the realizations is that of the product order:

$$x \lesssim y \quad \text{iff} \quad r_1(x) \leq r_1(y) \ \wedge \ r_2(x) \leq r_2(y).$$

We might ask whether we can "aggregate" the two rules into one. The answer is positive, given the conditions in the following definition.

**Definition 16** (Rule aggregation operation)**.** Consider a rulebook $\langle \mathcal{R}, \leq \rangle$ in which there are two rules $r_1, r_2 \in \mathcal{R}$ that are in the same equivalence class defined by $\leq$. Then it is allowed to "aggregate" the two rules into a new rule $r'$, defined by

$$r'(x) = \alpha(r_1(x), r_2(x)),$$

where $\alpha$ is an embedding of the product pre-order into $\mathbb{R}_+$.

In particular, allowed choices for $\alpha$ include linear combinations with positive coefficients ($\alpha(r_1, r_2) = a\, r_1 + b\, r_2$) and other functions that are strictly monotone in both arguments.

*4) Rule augmentation:* Adding a rule to a rulebook is a potentially destructive operation. In general, we can preserve the existing order only if the added rule is below every other.

**Definition 17** (Rule augmentation)**.** The operation of rule augmentation consists in adding to the rulebook $\mathcal{R}$ a rule $r'$ such that $r' < r$ for all $r \in \mathcal{R}$.

*5) Properties preserved by the three operations:* We will show that the three operations create a rulebook that is a refinement of the original rulebook, in the sense of Def. 18.

**Definition 18.** A rulebook $\langle \mathcal{R}_1, \leq_1 \rangle$ a *strict refinement* of $\langle \mathcal{R}_2, \leq_2 \rangle$ if its induced strict pre-order $<_2$ refines $<_1$.
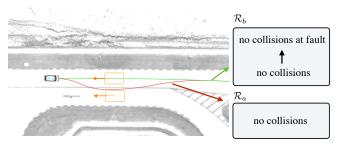
One can prove this theorem:



Fig. 7: Trajectories planned in the unavoidable collision scenario with different versions of the rulebooks. (See attached videos for experiment.) The orange rectangles are the traffic vehicles, moving towards the ego vehicle at the speed of 1.0 m/s. The red trajectory is chosen when collision at fault and collision caused by third-party are treated equally whereas the green trajectory is chosen when collision at fault is higher in the rulebooks hierarchy than the collision caused by third-party.

**Theorem 19.** *Applying one of the three operations (augmentation, refinement, aggregation) to a rulebook $\mathcal{R}_1$ creates a rulebook $\mathcal{R}_2$ that is a strict refinement of $\mathcal{R}_1$ in the sense of Def. 18.*

The proofs for these and ancillary results are in the appendix.

## V. Experiments

We show planning results for different rulebooks for the nuTonomy R&D platform (Renault Zoe). The experiments assume left-hand traffic (Singapore/UK regulations).

*1) Unavoidable collision:* This experiment illustrates unavoidable collision as described in Example 10. We set up the scenario (Fig. 7) such that the planner is led to believe that 2 vehicles instantaneously appear at approximately 12 m from the ego vehicle and slowly move towards the ego vehicle at 1.0 m/s, while the speed of the ego vehicle is 9.5 m/s. Fig. 7 shows the belief state when the vehicles first appear. We also limit the allowed deceleration to 3.5 m/s². It can be verified that collision is unavoidable under these conditions.

For any given trajectory $x$, we define the collision cost as

$$\mu(x) = v_{x,\text{col}}, \tag{1}$$

where $v_{x,\text{col}}$ is the expected longitudinal speed of the ego vehicle at collision, assuming that the ego vehicle applies the maximum deceleration from the current state.

First, consider the case where the collision cost (1) is applied to any collision. In this case, it is more preferable to swerve and hit the traffic vehicle in the opposite lane since the swerving trajectory gives the ego vehicle more distance to decelerate; hence, reducing the expected speed at collision.

Next, collision at fault $\mu_1$ is differentiated from collision caused by third-party $\mu_2$ with priority $\mu_2 < \mu_1$. The optimal trajectory in this case is to stay within lane and collide with the traffic vehicle that is moving against the direction of traffic.

*2) Clearance and lane keeping:* In this experiment, we demonstrate how different rulebooks in Example 11 lead to different behaviors when overtaking a stationary vehicle. The blockage cost $\beta$, lane keeping cost $\lambda$ and length $\alpha$ are defined as in Example 11, but we re-define the clearance cost as

$$\kappa(x) = \max(0, C_0 - l_x), \tag{2}$$

where $l_x$ is the minimum lateral distance between the stationary vehicle and trajectory $x$.

In particular, we consider two different rulebooks (Fig. 5):

$$\mathcal{R}_1 = \{\alpha < \lambda < \kappa < \beta\}, \quad \text{(clearance first)} \tag{3}$$
$$\mathcal{R}_2 = \{\alpha < \kappa < \lambda < \beta\}. \quad \text{(lane keeping first)} \tag{4}$$

The rulebook described in (4) corresponds to the case where satisfying the lane keeping rule is preferred over satisfying the clearance rule whereas the rulebook described in (3) corresponds to the case where satisfying the clearance rule is preferred over satisfying the lane keeping rule.

Fig. 8 shows the optimal paths found by the system in the two cases. With rulebook $\mathcal{R}_2$, the optimal trajectory is such

that the vehicle footprint remains within lane, leading to the violation of the clearance rule. In contrast, when rulebook $\mathcal{R}_1$ is applied, the trajectory crosses the lane boundary to give sufficient clearance from the stationary vehicle.
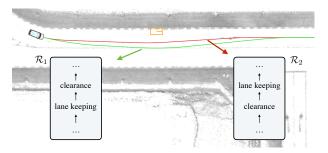


Fig. 8: Trajectories planned in the vehicle overtaking scenario with different rulebooks. (See attached videos for experiment.) The orange rectangle is the stationary vehicle. The red trajectory is when the rulebook (4) is used, whereas the green trajectory is when the rulebook (3) is used.

*3) Lane change near intersection:* Consider the scenario where the autonomous vehicle needs to perform a lane change in the vicinity of an intersection (Fig. 9). The vehicle needs to turn left at the intersection; therefore, it is required to be on the left lane before entering the intersection. However, there is a stationary vehicle that prevents it from completing the maneuver at an appropriate distance from the intersection.

For the simplicity of the presentation, we assume that any trajectory $x$ only crosses the lane boundary once at $\eta_x$. The lane change near intersection cost is then defined as $\zeta(x) = \max(0, D_{lc} - d_{\text{int}}(\eta_x))$, where $D_{lc}$ is a predefined threshold of the distance from intersection, beyond which changing lane is not penalized and for any pose $p$, $d_{\text{int}}(p)$ is the distance from $p$ to the closest intersection.

Additionally, we define the turning cost $\tau(x)$ as the $L_1$-norm of the heading difference between $x$ and the nominal trajectory associated with each lane. Consider the case where $\zeta$ and $\tau$ are in the same equivalence class and these rules are aggregated (16) as

$$r_{\zeta,\tau}(x) = \zeta(x) + c_\tau \tau(x), \qquad (5)$$

where $c_\tau > 0$ is a predefined constant. In this experiment, we consider the aggregated cost $r_{\zeta,\tau}$ and the blockage cost $\beta$ defined in Example 11 with priority $r_{\zeta,\tau} < \beta$. Fig. 9 shows how the choice of $c_\tau$ affects the optimal trajectory.
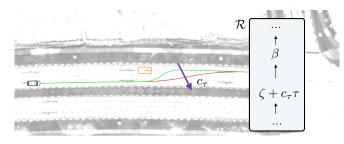


Fig. 9: Trajectories planned in the lane changing near intersection scenario. (See attached videos for experiment.) The orange rectangle is the stationary vehicle at pose $p_v$ with $d_{int}(p_v) < D_{lc}$. The green trajectory is the optimal trajectory for $c_\tau = 0$ whereas the red trajectory is the optimal trajectory for some $c_\tau > 0$.

## VI. DISCUSSION AND FUTURE WORK

We have shown by way of a few examples how the rulebooks approach allows easy and intuitive tuning of self-driving behavior. What is difficult to convey in a short paper is the ability of the formalism to scale up. In our production code at nuTonomy, corresponding to level 4 autonomy in a limited operating domain, our rulebooks have about 15 rules. For complete coverage of Massachusetts or Singapore rules, including rare corner cases (such as "do not scare farm animals"), we estimate about 200 rules, to be organized in about a dozen ordered priority groups (Fig. 10).

Except the extrema of safety at the top, all the other priorities among rule groups are somehow open for discussion. What we realized is that some of the rules and rules priorities, especially those that concern safety and liability, must *be part of nation-wide and global regulations* to be developed after an *informed public discourse*; it should not be up to engineers to choose these important aspects. The rulebooks formalism allows to have one such shared, high-level specification that gives minimal constraints to the behavior; then, the rest of the rules and priority choices can be considered "implementation details" that might change from manufacturer to manufacturer.
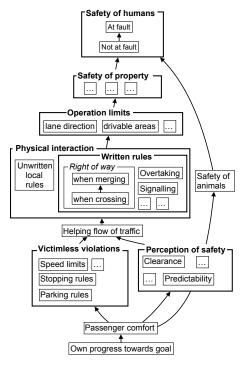


Fig. 10: Illustrative example of possible rule groups for an autonomous taxi in an urban driving scenario. At the top of the hierarchy there are rules that guarantee safety of humans; at the bottom, we have comfort constraints and progress goals. At the top, the rules are written analytically; at the bottom, some rules are learned from observed behavior. Rules at the bottom also tend to be platform- and implementation- specific. Except for human safety at the top, all other priorities among rule groups are open for discussion.

REFERENCES

[1] R. C. Arkin, "Governing Lethal Behavior: Embedding Ethics in a Hybrid Deliberative/Reactive Robot Architecture Part I: Motivation And Philosophy," *Proceedings of the 3rd international conference on Human robot interaction - HRI '08*, p. 121, jan 2008. [Online]. Available: http://portal.acm.org/citation.cfm?doid=1349822.1349839

[2] S. M. Thornton, S. Pan, S. M. Erlien, and J. C. Gerdes, "Incorporating ethical considerations into automated vehicle control," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 6, pp. 1429–1439, June 2017.

[3] E. Pires Bjørgen, S. Øvervatn Madsen, T. Skaar Bjørknes, F. Vonheim Heimsæter, R. Håvik, M. Linderud, P. Longberg, L. Dennis, and M. Slavkovik, "Cake, death, and trolleys: dilemmas as benchmarks of ethical decision-making," in *AAAI/ACM Conference on Artificial Intelligence, Ethics and Society*, New Orleans, USA, 2018, forthcoming.

[4] P. Foot, "The problem of abortion and the doctrine of double effect," *Oxford Review*, vol. 5, pp. 5–15, 1967.

[5] J. J. Thomson, "The trolley problem," vol. 94, no. 6, pp. 1395–1415, 1985. [Online]. Available: http://www.jstor.org/stable/796133

[6] "Moral machine (online)," 2016, http://moralmachine.mit.edu.

[7] G. Fainekos, H. Kress-Gazit, and G. Pappas, "Temporal logic motion planning for mobile robots," in *Proc. of IEEE International Conference on Robotics and Automation*, April 2005, pp. 2020–2025.

[8] M. Kloetzer and C. Belta, "A fully automated framework for control of linear systems from temporal logic specifications," *IEEE Transactions on Automatic Control*, vol. 53, no. 1, pp. 287–297, 2008.

[9] P. Tabuada and G. J. Pappas, "Linear time logic control of linear systems," *IEEE Transaction on Automatic Control*, vol. 51, no. 12, pp. 1862–1877, 2006.

[10] S. Karaman, R. G. Sanfelice, and E. Frazzoli, "Optimal control of mixed logical dynamical systems with linear temporal logic specifications," Dec. 2008, pp. 2117–2122.

[11] J. Liu, N. Ozay, U. Topcu, and R. M. Murray, "Synthesis of reactive switching protocols from temporal logic specifications," *IEEE Transactions on Automatic Control*, vol. 58, no. 7, pp. 1771–1785, 2013.

[12] M. Webster, M. Fisher, N. Cameron, and M. Jump, "Formal methods for the certification of autonomous unmanned aircraft systems," in *Proceedings of the 30th International Conference on Computer Safety, Reliability, and Security*, ser. SAFECOMP'11. Berlin, Heidelberg: Springer-Verlag, 2011, pp. 228–242. [Online]. Available: http://dl.acm.org/citation.cfm?id=2041619.2041644

[13] J. Morse, D. Araiza-Illan, K. Eder, J. Lawry, and A. Richards, "A fuzzy approach to qualification in design exploration for autonomous robots and systems," in *2017 IEEE International Conference on Fuzzy Systems, FUZZ-IEEE 2017, Naples, Italy, July 9-12, 2017*, 2017, pp. 1–6.

[14] I. Cizelj and C. Belta, "Negotiating the probabilistic satisfaction of temporal logic motion specifications," in *Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2013.

[15] H. Kress-Gazit, M. Lahijanian, and V. Raman, "Synthesis for robots: Guarantees and feedback for robot behavior," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, no. 1, pp. 211–236, 2018. [Online]. Available: https://doi.org/10.1146/annurev-control-060117-104838

[16] J. Van Den Hoven and G.-J. Lokhorst, "Deontic logic and computer-supported computer ethics," *Metaphilosophy*, vol. 33, no. 3, pp. 376–386. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1111/1467-9973.00233

[17] S. Levine and V. Koltun, "Continuous inverse optimal control with locally optimal examples," in *ICML '12: Proceedings of the 29th International Conference on Machine Learning*, 2012.

[18] L. I. R. Castro, P. Chaudhari, J. Tumova, S. Karaman, E. Frazzoli, and D. Rus, "Incremental sampling-based algorithm for minimum-violation motion planning," in *52nd IEEE Conference on Decision and Control*, Dec 2013, pp. 3217–3224.

[19] J. Tumova, G. C. Hall, S. Karaman, E. Frazzoli, and D. Rus, "Least-violating control strategy synthesis with safety rules," in *Proceedings of the 16th International Conference on Hybrid Systems: Computation and Control*, ser. HSCC '13. New York, NY, USA: ACM, 2013, pp. 1–10. [Online]. Available: http://doi.acm.org/10.1145/2461328.2461330

[20] J. Tumova, L. I. R. Castro, S. Karaman, E. Frazzoli, and D. Rus, "Minimum-violation LTL Planning with Conflicting Specifications," *American Control Conference*, jan 2013. [Online]. Available: http://arxiv.org/abs/1303.3679

[21] C. I. Vasile, J. Tumova, S. Karaman, C. Belta, and D. Rus, "Minimum-violation scLTL motion planning for mobility-on-demand," *Proceedings - IEEE International Conference on Robotics and Automation*, pp. 1481–1488, 2017.

[22] E. Frazzoli and K. Iagnemma, "US patent US9645577B1: Facilitating vehicle driving and self-driving."

[23] J. W. Polderman and J. C. Willems, *Introduction to Mathematical Systems Theory: A Behavioral Approach*. Berlin, Heidelberg: Springer-Verlag, 1998.

[24] R. C. Arkin, *Behavior-based Robotics*, 1st ed. Cambridge, MA, USA: MIT Press, 1998.

[25] M. J. Mataric and F. Michaud, "Behavior-based systems," in *Springer Handbook of Robotics*, 2008, pp. 891–909. [Online]. Available: https://doi.org/10.1007/978-3-540-30301-5_39

[26] V. Modugno, G. Neumann, E. Rueckert, G. Oriolo, J. Peters, and S. Ivaldi, "Learning soft task priorities for control of redundant robots," in *IEEE International Conference on Robotics and Automation (ICRA 2016)*, Stockholm, Sweden, May 2016. [Online]. Available: https://hal.archives-ouvertes.fr/hal-01273409

[27] J. Silvério, S. Calinon, L. D. Rozo, and D. G. Caldwell, "Learning competing constraints and task priorities from demonstrations of bimanual skills," *CoRR*, vol. abs/1707.06791, 2017. [Online]. Available: http://arxiv.org/abs/1707.06791

[28] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "On a formal model of safe and scalable self-driving cars," 08 2017.

[29] A. Aswani, H. Gonzalez, S. S. Sastry, and C. Tomlin, "Provably safe and robust learning-based model predictive control," *Automatica*, vol. 49, no. 5, pp. 1216–1226, jan 2013. [Online]. Available: http://linkinghub.elsevier.com/retrieve/pii/S0005109813000678papers3://publication/doi/10.1016/j.automatica.2013.02.003

# APPENDIX I
## ORDERS AND PREORDERS

This appendix recalls some standard definitions used in the development of the rulebooks formalism.

**Definition 20.** A *preorder* on a set $Z$ is a reflexive transitive binary relation $\lesssim$, i.e., a binary relation such that for all $z \in Z$ one has $z \lesssim z$ and for all $x, y, z \in Z$

$$x \lesssim y \text{ and } y \lesssim z \implies x \lesssim z.$$

Given a preorder $(Z, \lesssim)$ and $y, z \in Z$, the notation $y < z$ is shorthand for $y \lesssim z$ and $z \not\lesssim y$. A preorder is said to be *total* if additionally for any $x, y \in Z$ either $x \lesssim y$ or $y \lesssim x$.

With any preorder one associates an equivalence relation: elements $x, y \in Z$ are equivalent (denoted as $x \sim y$) whenever $x \lesssim y$ and $y \lesssim x$. If this equivalence relation is trivial (i.e., $x \sim y$ if and only if $x = y$), then we say that $\lesssim$ is a *partial order* on $Z$. We use $x \leq y$ to denote ($x < y$ or $x = y$). In particular, one always has $\leq \subseteq \lesssim$, and a preorder is an order if and only if $\leq = \lesssim$.

**Definition 21.** An *embedding* between preorders $Z_1$ and $Z_2$ is a map $\phi : Z_1 \to Z_2$ such that for all $x, y \in Z_1$ one has

$$x \lesssim y \implies \phi(x) \lesssim \phi(y) \text{ and } x < y \implies \phi(x) < \phi(y).$$

Note that for any embedding $\phi$, equivalence $x \sim y$ implies $\phi(x) \sim \phi(y)$.

**Definition 22.** Given a set $Z$ and two preorders $\lesssim_1$, $\lesssim_2$ on it, we say that $\lesssim_2$ *refines* $\lesssim_1$ if the identity map $\text{id} : Z \to Z$ is an embedding from $(Z, \lesssim_1)$ onto $(Z, \lesssim_2)$.

# APPENDIX II
## OPERATIONS ON RULEBOOKS

As explained in Section 2, a rulebook induces a partial preorder $\lesssim$ on realizations. We are interested in operations on the rulebooks that preserve existing relations, but may possibly introduce new comparisons between realizations, i.e., the preorder on realizations may be refined.

More formally, our goal is to find conditions on a map $\phi : \mathcal{R}_1 \to \mathcal{R}_2$ between two rulebooks that guarantee that $\lesssim_2$ is a refinement of $\lesssim_1$. To motivate concepts that will follow, let us begin with the simplest case of a rulebook $\mathcal{R}_2 = \{u\}$ consisting of a single rule. Let $\{r_1, \ldots, r_n\} = \mathcal{R}_1$ be the rules in the domain, and the map $\phi$ therefore collapses all $r_i$ onto $u$: $\phi(r_i) = u$ for all $i$. At the moment we do not impose any assumptions on $r_i$ — some of them may be comparable, some are equivalent or independent. The question then becomes when the (total) preorder imposed by $u$ on $\Xi$ is a refinement of the (partial) preorder given by $\{r_1, \ldots, r_n\}$. Recall that according to the definition of a refinement that amounts to:

$$\forall x, y \in \Xi \quad (x \lesssim_1 y \implies x \lesssim_2 y) \quad \text{and}$$
$$(x <_1 y \implies x <_2 y). \quad (6)$$

## A. Aggregative maps

The first observation is that equation (6) necessarily implies that the value $u(x)$ depends only on the values $r_1(x), \ldots, r_n(x)$ and not on the realization $x$ itself. Indeed, if $x$ and $y$ are two realizations such that $r_i(x) = r_i(y)$ for all $i$, yet $u(x) \neq u(y)$ (say, $u(x) < u(y)$ for definiteness), then $y \lesssim_1 x$, but $y \not\lesssim_2 x$, contradicting equation (6).

One may view $\mathcal{R}_1$ as providing a map from all realizations to $(\mathbb{R}^+)^n$ via

$$x \mapsto (r_1(x), \ldots, r_n(x)),$$

and similarly $\mathcal{R}_2$, consisting just of a single rule, can be identified with a map $u : \Xi \to \mathbb{R}^+$. The observation above can then be reinterpreted to say that there exists a map $\alpha : (\mathbb{R}^+)^n \to \mathbb{R}^+$ making the following diagram commutative:

$$\forall x \in \Xi \quad \alpha(r_1(x), \ldots, r_n(x)) = u(x).$$



Fig. 11: Factorization of $u : \Xi \to \mathbb{R}^+$

What can be said about the map $\alpha$ itself? The set $(\mathbb{R}^+)^n$ has a natural partial order on it, called the product order, where given $\vec{a}, \vec{b} \in (\mathbb{R}^+)^n$,

$$\vec{a} = (a_1, \ldots, a_n), \quad \vec{b} = (b_1, \ldots, b_n),$$

one denotes $\vec{a} \leq \vec{b}$ whenever $a_i \leq b_i$ for all $i$. Note that if $x, y \in \Xi$ are two realizations such that

$$(r_1(x), \ldots, r_n(x)) \leq (r_1(y), \ldots, r_n(y)),$$

then necessarily $x \lesssim_1 y$ and in view of equation (6) we therefore need to have $u(x) \leq u(y)$. If moreover $r_i(x) < r_i(y)$ for at least some $i$, then $x <_1 y$, and so $u(x) < u(y)$ must be the case. This observation can be summarized as follows: If $S \subseteq (\mathcal{R}^+)^n$, $T \subseteq \mathbb{R}^+$ are the sets

$$S = \{(r_1(x), \ldots, r_n(x)) : x \in \Xi\}, \quad T = \{u(x) : x \in \Xi\},$$

then $\alpha : S \to T$ is an embedding of partial orders. This brings us to the following

**Definition 23** (Aggregative map)**.** We say that it is *admissible to collapse rules* $r_1, \ldots, r_n$ to a rule $u$ if, in the notation above, the map $\alpha$, that makes Figure 11 commutative, exists, and $\alpha : S \to T$ is an embedding of partial orders.

A map between rulebooks $\phi : \mathcal{R} \to \mathcal{R}'$ is said to be *aggregative* if for all $u \in \phi(\mathcal{R}')$ it is admissible to collapse rules $\phi^{-1}(u)$ onto $u$.

The following lemma shows that *surjective* aggregative maps can be composed yielding another surjective aggregative map.

**Lemma 24.** *Composition of surjective aggregative maps is aggregative.*

*Proof.* Let $\mathcal{R}, \mathcal{R}', \mathcal{R}''$ be rulebooks, let $\phi_1 : \mathcal{R} \to \mathcal{R}'$, $\phi_2 : \mathcal{R}' \to \mathcal{R}''$ be aggregative maps and let $\psi : \mathcal{R} \to \mathcal{R}''$ be the composition of the two, $\psi = \phi_2 \circ \phi_1$. We need to show that $\psi$ is aggregative and to this end pick some $w \in \mathcal{R}''$. Let $u_1, \ldots, u_n$ be the rules in $\phi_2^{-1}(w)$, and for each $1 \leq i \leq n$ let $r_1^i, \ldots, r_{m_i}^i$ be the rules enumerating $\phi_1^{-1}(u_i)$ (see Figure 12).



Fig. 12: Structure of the preimage $\psi^{-1}(w)$ of $w$

The first observation is that for any realization $x \in \Xi$, the value $w(x)$ depends only on the numbers $r_j^i(x)$. Indeed, since $\phi_1$ is aggregative, for each $i$ the value $u_i(x)$ depends only on $(r_j^i(x))_{1 \leq j \leq m_i}$, and by a similar token $w(x)$ is uniquely reconstructible from $u_i(x)$. More precisely, suppose $\alpha_i : (\mathbb{R}^+)^{m_i} \to \mathbb{R}^+$ are the maps witnessing that $\phi_1$ is aggregative, and $\beta : (\mathbb{R}^+)^n \to \mathbb{R}^+$ is the corresponding map for $\phi_2$. If $t = \sum_{i=1}^n m_i$, then the map $\gamma : (\mathbb{R}^+)^t \to \mathbb{R}^+$ given by

$$
\begin{aligned}
\gamma(c_1, &\ldots, c_t) \\
&= \beta(\alpha_1(c_1, \ldots, c_{m_1}), \alpha_2(c_{m_1+1}, \ldots, c_{m_1+m_2}), \ldots, \\
&\qquad\qquad \alpha_n(c_{t-m_n+1}, \ldots, c_t)),
\end{aligned}
$$

satisfies

$$
\begin{aligned}
\gamma(r_1^1(x), &\ldots, r_{m_1}^1(x), \ldots, r_1^n(x), \ldots, r_{m_n}^n(x)) \\
&= \beta(u_1(x), \ldots, u_n(x)) = w(x)
\end{aligned}
$$

for all $x \in \Xi$.

We need to show that $\gamma$ is an embedding of partial orders, and to this end let $x, y \in \Xi$ be realizations such that $r_j^i(x) \leq r_j^i(y)$ for all $i, j$. We need to show that $w(x) \leq w(y)$. Note that since $\alpha_i$'s are embeddings,

$$
\begin{aligned}
u_i(x) = \alpha_i(r_1^i(x), \ldots, r_{m_i}^i(x)) &\leq \\
\alpha_i(r_1^i(y), \ldots, r_{m_i}^i(y)) &= u_i(y).
\end{aligned}
$$

Also, since $\beta$ is an embedding, this implies that

$$
\begin{aligned}
w(x) = \beta(u_1(x), \ldots, u_n(x)) &\leq \\
\beta(u_1(y), \ldots, u_n(y)) &= w(y),
\end{aligned}
$$

and hence $w(x) \leq w(y)$ as claimed.

Finally, if moreover $r_j^i(x) < r_j^i(y)$ for some $i, j$, then $u_i(x) < u_i(y)$, and therefore also $w(x) < w(y)$. Thus $\gamma$ is

an embedding of partial orders, and therefore $\psi$ is aggregative. $\qquad\square$

We are now ready to introduce the key notion of an embedding between rulebooks.

**Definition 25** (Rulebook embedding)**.** An *embedding* between rulebooks is an aggregative map $\phi : \mathcal{R} \to \mathcal{R}'$ that is also an embedding between $\mathcal{R}$ and $\mathcal{R}'$ as partially preordered sets.

**Lemma 26.** *Let $\mathcal{R}_1$ and $\mathcal{R}_2$ be rulebooks, and let $\phi : \mathcal{R}_1 \to \mathcal{R}_2$ be a surjective embedding. Let $x, y \in \Xi$ be realizations such that $x \lesssim_1 y$. If $r \in \mathcal{R}_1$ is such that $r(y) \neq r(x)$, then there exists $u' \in \mathcal{R}_2$, such that*

$$
u' \geq \phi(r) \text{ and } u'(x) < u'(y).
$$

*Proof.* If $r(y) < r(x)$, then there exists $r_1 > r$ such that $r_1(x) < r_1(y)$. If $r(x) < r(y)$ to begin with, then we set $r_1 = r$. In either case we have $r_1 \geq r$ and $r_1(x) < r_1(y)$. Set $u_1 = \phi(r_1)$. We are done if $u_1(x) < u_1(y)$. Otherwise, let $\{r_1^1, \ldots, r_{m_1}^1\} = \phi^{-1}(u_1)$ be the preimage of $u_1$ (note that $r_1$ is one of these elements)

Since $\phi$ is aggregative and since $r_1(x) < r_1(y)$, there has to be some $1 \leq i \leq m_1$ such that $r_i^1(y) < r_i^1(x)$. Indeed, if $r_j^1(x) \leq r_j^1(y)$ for all $j$, then

$$
(r_1^1(x), \ldots, r_{m_1}^1(x)) < (r_1^1(y), \ldots, r_{m_1}^1(y))
$$

in the product order. Hence, $\phi$ being aggregative implies $u_1(x) < u_1(y)$ contradicting our earlier assumption. Thus $r_i^1(y) < r_i^1(x)$ for some $i$.

In view of $x \lesssim_1 y$, there exist $r_2 > r_i^1$ such that $r_2(x) < r_2(y)$. Set $u_2 = \phi(r_2)$. Note that

$$
r_i^1 < r_2 \implies \phi(r_i^1) < \phi(r_2) \iff u_1 < u_2.
$$

We are done if $u_2(x) < u_2(y)$.

Suppose that $u_2(x) \geq u_2(y)$ and let

$$
\{r_1^2, \ldots, r_{m_2}^2\} = \phi^{-1}(u_2).
$$

By the same argument as above, there must exist some $1 \leq i \leq m_2$ such that $r_i^2(x) > r_i^2(y)$. In view of $x \lesssim_1 y$, there is $r_3 > r_i^2$ such that $r_3(x) < r_3(y)$. Set $u_3 = \phi(r_3)$.
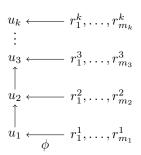


Fig. 13: Construction of the chain

The process continues, and builds a sequence of rules $u_1 < u_2 < \cdots < u_k$ as in Figure 13. By finiteness of the rulebook, the chain has to stop at some point, which is

possible only if $u_k(x) < u_k(y)$. Since $u_k > u_1 \geq \phi(r)$, the lemma follows. □

**Theorem 27.** *Let $\mathcal{R}_1$ and $\mathcal{R}_2$ be rulebooks. If there exists a surjective embedding $\phi : \mathcal{R}_1 \to \mathcal{R}_2$ between the two, then $\lesssim_2$ refines $\lesssim_1$.*

*Proof.* Suppose $x, y \in \Xi$ are such that $x \lesssim_1 y$, we show that $x \lesssim_2 y$. Pick some $u \in \mathcal{R}_2$ such that $u(y) < u(x)$. Let $\{r_1, \ldots, r_m\} = \phi^{-1}(u)$ be the preimage. Note that $r_i(x) > r_i(y)$ for some $i$ (for otherwise $u(x) \leq u(y)$, because $\phi$ is aggregative), hence Lemma 26 applies and produces some $u' \geq \phi(r_i) = u$ such that $u'(x) < u'(y)$.

It remains to show that $x <_1 y$ implies $x <_2 y$. Since $x \lesssim_2 y$ has already been shown, it is enough to show that $u(x) \neq u(y)$ for some $u \in \mathcal{R}_2$. Pick some $r \in \mathcal{R}_1$ such that $r(x) < r(y)$. Lemma 26 produces $u \geq \phi(r)$ such that $u(x) < u(y)$. □

Two examples of surjective embeddings are Priority Refinements (Def. 14) and Rule Aggregation (Def. 16).

## APPENDIX III
### ADDING NEW RULES

There is one important operation that is missing from the picture — addition of new rules. Surjective embeddings of rulebooks let us impose new relations between existing rules, as well as to aggregate several rules into one. What if one would like to add a new rule that does not bear any direct relation to existing ones?

Generally, this is a very destructive operation in the sense that it can dramatically change the preorder imposed on realizations. Perhaps the most extreme example is when to a rulebook $\mathcal{R}_1$ a new rule $r$ is added that is declared to be of the highest importance: $u < r$ for all $u \in \mathcal{R}_1$. Let $\mathcal{R}_2$ denote the resulting rulebook $\mathcal{R}_1 \cup \{r\}$, and note that if $x, y$ are two realizations such that $r(x) < r(y)$ then necessarily $x <_2 y$ regardless of how they were related in the preorder induced by $\mathcal{R}_1$.

A similar but slightly more general case is when a new rule is added "in the middle" of $\mathcal{R}_1$. More formally, suppose that $\mathcal{R}_2 = \mathcal{R}_1 \cup \{r\}$, where the new rule $r$ satisfies $u < r$ for some $u \in \mathcal{R}_1$. If $x, y \in \Xi$ are two realizations such that $u'(x) = u'(y)$ for all $u' \in \mathcal{R}_1 \setminus \{u\}$ and $u(x) < u(y)$, then $x <_1 y$. If, however, $r(y) < r(x)$, then $y <_2 x$, and the order between the realizations is reverted.

Unless one makes some additional assumptions on the set of realizations $\Xi$, the example above shows that adding a rule above an existing one can easily change the preorder on realizations. However, in order to get some meaningful preservation of the preorder, it is not enough to assume that the newly added rules are not above any of the existing ones. Consider the simplest case, when $\mathcal{R}_1$ consists of a single rule $\{u\}$, and $\mathcal{R}_2 = \{r, u\}$ adds a rule that is incomparable with $r$. If $x$ and $y$ are two realizations such that $u(x) < u(y)$, then necessarily $x <_1 y$. However, if $r(y) < r(x)$, then $x$ and $y$ are incomparable relative to $\lesssim_2$. When the two rules are further aggregated as described in the previous appendix, all relations between $x$ and $y$ become possible.

The example above can be modified slightly by considering a rulebook $\mathcal{R}_1 = \{u, u'\}$, $u < u'$, and adding the rule $r$ such that $r < u'$, but $r$ and $u$ are incomparable. The same analysis as above now applies to a pair of realizations such that $u'(x) = u'(y)$. In particular, for the relation $\lesssim$ to be broken by adding a new rule, one does not have to add a completely independent rule, it is enough to have some rules in $\mathcal{R}_1$ that are not comparable to $r$.

We are left with only one option — add new rules below all of the existing ones. However, even this operation does not result in the refinement of the $\lesssim$ order on realizations. Indeed, the simplest case, is when $\mathcal{R}_1 = \{u\}$ and $\mathcal{R}_2 = \{u, r\}, r < u$. If $x, y$ are two equivalent realizations, then necessarily $x \lesssim_1 y$ and $y \lesssim_1 x$. However, if the new rule $r$ differentiates between the realizations, $r(x) \neq r(y)$, then one of $x \lesssim_2 y$, $y \lesssim_2 x$ is false.

We conclude that in general we cannot guarantee that the relation $\lesssim$ has been refined if any new rules were added. The last example in the list above is, nonetheless, different from others. It turns out that the only problem that can occur, when new rules are added below existing ones, is that equivalent realizations are no longer equivalent in the enlarged rulebook. Thus, while the preorder $\lesssim$ may not be refined, its strict counterpart $<$ is preserved by such an operation.

**Definition 28.** An embedding of rulebooks $\phi : \mathcal{R}_1 \to \mathcal{R}_2$ is said to be *dominant*, if $u < r$ for all $u \in \mathcal{R}_2 \setminus \phi(\mathcal{R}_1)$ and $r \in \phi(\mathcal{R}_1)$.

**Theorem 29.** *Let $\phi : \mathcal{R}_1 \to \mathcal{R}_2$ be a dominant embedding of rulebooks. If $x, y$ are realizations such that $x <_1 y$, then also $x <_2 y$.*

*Proof.* Consider $\phi(\mathcal{R}_1)$ as a rulebook, and note that the map $\phi : \mathcal{R}_1 \to \phi(\mathcal{R}_1)$ is automatically surjective. Theorem 27 applies, and shows that $x < y$ relative to $\phi(\mathcal{R}_1)$ as well. This allows us to assume without loss of generality that $\mathcal{R}_1 \subseteq \mathcal{R}_2$, and the map $\phi$ is the identity map. Since $\phi$ is assumed to be dominant, it means that $u < r$ for all $r \in \mathcal{R}_1$ and $u \in \mathcal{R}_2 \setminus \mathcal{R}_1$.

Suppose $x, y$ are two realizations such that $x <_1 y$, and let $r \in \mathcal{R}_2$ be such that $r(y) < r(x)$. We need to show that there exists some rule $r' > r$ such that $r'(x) < r'(y)$. Indeed, if $r \in \mathcal{R}_1$, then such a rule $r'$ must exist simply because $x <_1 y$ by assumption. So, let us assume that $r \in \mathcal{R}_2 \setminus \mathcal{R}_1$. Since $x <_1 y$ there mush be at least one rule $u \in \mathcal{R}_1$ such that $u(x) < u(y)$. Since $u > r$, the theorem follows. □

Rule Augmentation as described in Def. 17 is an example of a dominant embedding.