

Toward the Engineering of Virtuous Machines

Naveen Sundar Govindarajulu
Rensselaer AI & Reasoning Lab
Rensselaer Polytechnic Institute, (RPI)
Troy, NY, USA
naveensundarg@gmail.com

Rikhiya Ghosh
Rensselaer AI & Reasoning Lab
Rensselaer Polytechnic Institute, (RPI)
Troy, NY, USA
rikrixa@gmail.com

Selmer Bringsjord
Rensselaer AI & Reasoning Lab
Rensselaer Polytechnic Institute, (RPI)
Troy, NY, USA
Selmer.Bringsjord@gmail.com

Vasanth Sarathy
Human Robot Interaction Laboratory
Tufts University
Medford, MA, USA
vasanth.sarathy@tufts.edu

ABSTRACT

While various traditions under the ‘virtue ethics’ umbrella have been studied extensively and advocated by ethicists, it has not been clear that there exists a version of virtue ethics rigorous enough to be a target for machine ethics (which we take to include the engineering of an ethical sensibility in a machine or robot itself, not only the study of ethics in the humans who might create artificial agents). We begin to address this by presenting an embryonic formalization of a key part of any virtue-ethics theory: namely, the learning of virtue by a focus on exemplars of moral virtue. Our work is based in part on a computational formal logic previously used to formally model other ethical theories and principles therein, and to implement these models in artificial agents.

CCS CONCEPTS

• **Theory of computation** → **Proof theory; Modal and temporal logics; Logic and verification; Higher order logic.**

KEYWORDS

logic; virtue ethics; deontic cognitive event calculus; virtuous robots; verification

ACM Reference Format:

Naveen Sundar Govindarajulu, Selmer Bringsjord, Rikhiya Ghosh, and Vasanth Sarathy. 2019. Toward the Engineering of Virtuous Machines. In *AAAI/ACM Conference on AI, Ethics, and Society (AI/ES '19)*, January 27–28, 2019, Honolulu, HI, USA. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3306618.3314256>

INTRODUCTION

What is virtue ethics? One way of summarizing virtue ethics is to contrast it with the two main families of ethical theories: **deontological ethics** (\mathcal{D}) and **consequentialism** (\mathcal{C}). Ethical theories in the family \mathcal{C} that are utilitarian in nature hold that actions are

morally evaluated based on their **total utility** (or total *disutility*) to everyone involved. The best action is the action that has the highest total utility. In stark contrast, ethical theories in \mathcal{D} emphasize **inviolable principles**, and reasoning from those principles to whether actions are obligatory, permissible, neutral, etc.¹ In a departure from both \mathcal{D} and \mathcal{C} , ethical theories in the virtue-ethics family \mathcal{V} are overall distinguished by the principle that the best action in a situation, morally speaking, is the one that a **virtuous person** would perform [3]. A virtuous person is defined as a person who has learned and internalized a set of habits or traits termed **virtuous**. For a virtuous person, virtuous acts become second-nature, and hence are performed in many different situations, through time.

While extensive formal and rigorous modeling under the umbrellas of \mathcal{D} and \mathcal{C} has been carried out, there has been little such work devoted to formalizing and mechanizing \mathcal{V} . Note that unlike \mathcal{D} and \mathcal{C} , it is not entirely straightforward how one could translate the concepts and principles in \mathcal{V} into a form that is precise enough to be realized in machines. Proponents of \mathcal{V} might claim that it is not feasible to do so given \mathcal{V} 's emphasis on persons and traits, rather than individual actions or consequences. From the perspective of machine ethics, this is not satisfactory. If \mathcal{V} is to be on equal footing with \mathcal{D} and \mathcal{C} for the purpose of building morally competent machines, AI researchers need to start formalizing parts of virtue ethics, and to then implement such formalization in computation.

We present one such formalization herein; one that uses learning and is based on a virtue-ethics theory presented by Zagzebski [31]. The formalization is presented courtesy of an expressive computational logic that has been used to model principles in both \mathcal{C} and \mathcal{D} [e.g. [12, 14]].² The formalization answers, abstractly, the following two questions:

Questions

- (Q₁) When can we say an agent is virtuous?
- (Q₂) What is a virtue?

The plan for the paper is as follows. First, we briefly consider why virtuous machines might be useful, and then we briefly cover related work that can be regarded formalization of virtue ethics. Next, we

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
AI/ES '19, January 27–28, 2019, Honolulu, HI, USA
© 2019 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-6324-2/19/01.
<https://doi.org/10.1145/3306618.3314256>

¹Both the families \mathcal{C} and \mathcal{D} are crisply explained as being in conformity with what we say here in e.g. [9].

²See [6] for an introduction to the logicist methodology for building ethical machines.

present an overview of virtue ethics itself, and specifically show that an emphasis on moral exemplars makes good sense for any attempt to engineer a virtuous machine. We next present one version of virtue ethics, \mathcal{V}_z (Zagzebski’s version of virtue ethics), which we seek to formalize fully. Then, our calculus and the formalization itself (\mathcal{V}_z^f) are presented. We conclude by discussing future work and remaining challenges.

WHY VIRTUOUS ROBOTS?

Note that we do not advocate that machine ethicists pursue virtue ethics over other families of ethical theories. Our goal in the present paper is merely to formalize one version of virtue ethics within the family \mathcal{V} . That said, why might virtue ethics be preferred over consequentialism or deontological ethics for building morally competent machines? To partially answer this question, we take a short digression into a series of conditions laid out by Alfano, and characterized as identifying the core of virtue ethics:

Hard Core of Virtue Ethics (partially quoting [2])

- (2) **stability** If someone possesses a virtue at time t_1 , then *ceteris paribus* she will possess that virtue at a later time t_2 .
- (3) **consistency** If someone possesses a virtue sensitive to reason r , then *ceteris paribus* she will respond to r in most contexts.
- (7) **explanatory power** If someone possesses a virtue, then reference to that virtue will sometimes help to explain her behavior.
- (8) **predictive power** If someone possesses a high-fidelity virtue, then reference to that virtue will enable nearly certain predictions of her behavior; if someone possesses a low-fidelity virtue, then reference to that virtue will enable weak predictions of her behavior.

Particularly, we feel that if the conditions of **stability**, **consistency**, **explanatory power**, and **predictive power** hold, then virtuous agents or robots might be easier for humans to understand and interact with (compared to consequentialist or deontological agents or robots). This is but our initial motivation; we now present an overview of virtue ethics, in order to show that our focus specifically on learning of virtuous behavior from moral exemplars is advisable.

SURVEYING VIRTUE ETHICS

See [28] for a general introduction to the field of moral robots. We begin our survey by reporting that Hurka [16] presents an ingenious formal account involving a recursive notion of goodness and badness. The account starts with a given set of primitive good and bad states-of-affairs. Virtues are then defined as love of good states-of-affairs or hatred of bad states-of-affairs. Vice is defined as love of bad states-of-affairs or hatred of good states-of-affairs. Virtues and vices are then themselves taken to be good and bad states-of-affairs, resulting in a recursive definition (see Figure 2) that is attractive to AI researchers and computer scientists. But despite this, and despite our sense that the main problems with Hurka’s account are rectifiable [15], we feel that Hurka’s definition

might not capture central aspects of virtue [19]. More problematic is that it must be shown that Hurka’s account is different from rigorous and formal accounts of \mathcal{C} , which after all are themselves invariably based upon good and bad states-of-affairs. Moreover, it is not clear to us how Hurka’s account is amenable to automation. Therefore, we now proceed to step back and survey the overarching family \mathcal{V} of virtue ethics, to specifically pave a more promising AI road: viz. a focus on moral exemplars.

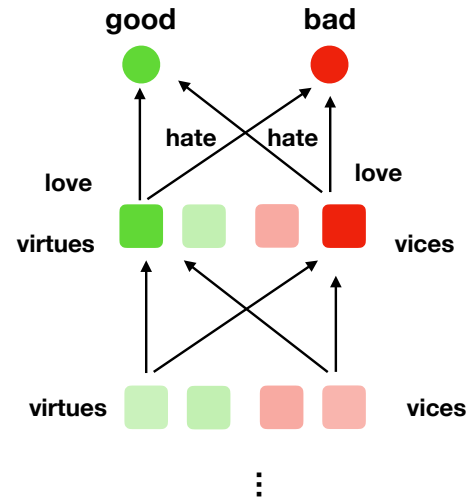


Figure 1: Hurka’s Account Virtues (vices) are defined recursively as love of good (bad) states-of-affairs or hate (love) of bad states of affairs.

Virtue Ethics: Overview to Exemplarism

The core concepts of consequentialist ethical theories (i.e. of members of \mathcal{C}), at least certainly in the particular such theory known as *utilitarianism*, are doubtless at minimum relatively familiar to most of our readers. For instance, most in our audience will know that utilitarianism’s core tenet is that actions are obligatory just in case they have the consequence of maximizing happiness, and are forbidden exactly when they fail to so maximize. A parallel state-of-affairs holds for at least basic knowledge of deontological ethical theories (= family \mathcal{D}): most readers have for instance some familiarity with Kant’s moral system in family \mathcal{D} , and specifically with his famous “categorical imperative,” which, paraphrasing, says that, unconditionally, one must always act in such a way that this behavior could be universalized.³ In addition, generally people are familiar with the core tenet of divine-command ethical theories (i.e. of members of \mathcal{DC}), which is (approximately) that actions are obligatory for humans if and only if God commands that these actions be performed (a particular member of \mathcal{D} is specified in [25]; another member is used for AI/machine ethics in [8]). However, in

³This imperative is first set out in — as it’s known in abbreviation — *Groundwork*; see [18]. It’s generally thought by ethicists, and this may be convenient for machine/AI ethics, that Kant had in mind essentially a decision procedure to follow in the attempt to behave in an ethically correct manner. For a lucid and laconic overview of this point, see [17]; and cf. [24].

our experience the epistemic situation is radically different when it comes to the family of ethical theories *virtue ethics* (= \mathcal{V}). For while it's true that generally educated people can be assumed to be acquainted with the pre-analytic concept of virtue, and with many things long deemed to be virtues (e.g. bravery), an understanding of virtue ethics at the level of ethical *theory* cannot be assumed. We therefore now provide a rapid (and admittedly cursory) synopsis of \mathcal{V} , by drawing from [29], and to some degree from [3]. It will be seen that \mathcal{V} makes central use of exemplars, and of learning and development that revolves around them. Hence we shall arrive at a convenient entry point for our AI work devoted to trying to design and build a virtuous machine.

Obviously we cannot in the span of the space we have at hand do full justice to the book-length treatment of \mathcal{V} that is [29]. But we can quickly establish that our technical work, in its focus on the cultivation of virtue for a machine via learning from exemplars, is not merely based on a single, idiosyncratic member of \mathcal{V} , and on one peripheral aspect of this member. On the contrary, study of the work of Vallor and other scholars concerned with a characterization of the family \mathcal{V} confirms that our exploitation specifically of Zagzebski's [31] focus, from the standpoint of the field of ethics itself, is a worthy point of entry for AI researchers.

To begin, Vallor, drawing on and slightly adapting Van Norden's [30], sets out a quartet of commonalities that at least seem to be true of all members of \mathcal{V} , and the second one is: "A conception of moral virtues as cultivated states of character, manifested by those exemplary persons who have come closest to achieving the highest human good" (§5, §2.2).⁴ But given our specific efforts toward engineering a virtuous machine, it is important to note that Vallor specifically informs us about the key concepts of exemplars in the particular members of the \mathcal{V} family; to pick just one of many available places, she writes:

Buddhism's resonances with other classical virtue traditions do not end here. As with the central role granted by Confucian and Aristotelian ethics to 'exemplary persons' (the *junzi* and *phronimoi* respectively), *bodhisattvas* (persons actively seeking enlightenment) generally receive direction to or assistance on the path of self-cultivation from the community of exemplary persons to which they have access. In Buddhism this is the monastic community and lay members of the *Sangha* . . . [¶5, §2.1.3, [29]]

We said above that we would also draw, albeit briefly, from a second treatment of \mathcal{V} , viz. [3], in order to pave the way into our AI-specific, exemplar-based technical work. About this second treatment we report only that it is one based squarely on a "range of development" (¶3, §Right Action in Ch. 3), where the agent (a human in her case) gradually develops into a truly virtuous person, beginning with unreflective adoption of direct instruction, through a final phase in which "actions are based on understanding gained through experience and reflection" (ibid.). Moreover, Annas explicitly welcomes the analogy between an agent's becoming virtuous,

⁴In her book, Vallor gives her own more detailed and technologically relevant list of seven core elements that can be viewed as common to all members of \mathcal{V} (or two what she refers to as "traditions" within virtue ethics). We do not have the space to discuss this list, and show that it fits nicely with our technical work's emphasis on exemplars and learning therefrom.

and an agent's becoming, say, an excellent tennis-player or pianist. The idea behind the similarity is that "two things are united: the *need to learn* and the *drive to aspire* (emphasis hers; ¶4 Ch. 3). In addition, following Aristotle on \mathcal{V} (e.g. see [4] 1103), no one can become a master tennis-player or pianist without, specifically, playing tennis/the piano with an eye to the mastery of great exemplars in these two domains.

In order to now turn to specific AI work devoted to engineering a virtuous machine, we move from completed consideration of the general foundation of \mathcal{V} , and its now-confirmed essential use of moral exemplars, to a specific use of such exemplars that appears ripe for mechanization.

EXEMPLARIST VIRTUE THEORY

Exemplarist virtue theory (\mathcal{V}_z) builds on the **direct reference theory** (DRT) of semantics. Briefly, in DRT, given a word or term w , its meaning $\mu(w)$ is determined by what the word picks out, say p , and not by some definition d . For example, for a person to use the word "water," in a correct manner, that person neither needs to possess a definition of water nor needs to understand all the physical properties of water. The person simply needs to know which entity the word "water" denotes in common usage.

In \mathcal{V}_z , persons understand moral terms, such as "honesty," in a similar manner. That is, moral terms are understood by persons through direct references instantiated in **exemplars**. Persons identify moral exemplars through the emotion of **admiration**. Fittingly, the emotion of admiration plays a foundational role in this theory (as does contempt). \mathcal{V}_z posits a process very similar to scientific or empirical investigation. Exemplars are first identified and their traits are studied; then they are continuously further studied to better understand their traits, qualities, etc. The status of an individual as an exemplar can change over time. Below is an informal version that we seek to formalize:

Informal Version \mathcal{V}_z

- I₁ Agent or person a perceives a person b perform an action α . This observation causes the emotion of admiration in a .
- I₂ a then studies b and seeks to learn what traits (habits or dispositions) b has.

THE GOAL

From the above presentation of \mathcal{V}_z , we can glean the following distilled requirements that should be present in any formalization.

\mathcal{V}_z^f Formalization Components

- (R₁) A formalization of emotions, particularly admiration.
- (R₂) A representation of traits.
- (R₃) A process of learning traits (and not just simple individual actions) from a small number of observations.

BUILDING THE FORMALIZATION

For fleshing out the above requirements and formalizing \mathcal{V}_z , we use the **deontic cognitive event calculus** (*DC $\mathcal{E}\mathcal{C}$*), a computational formal logic. This logic was used previously in [12, 14] to automate versions of the Doctrine of Double Effect (*DDE*), an ethical principle with deontological, consequentialist (and, historically, divine-command) components. *DC $\mathcal{E}\mathcal{C}$* has also been used to formalize *akrasia* (the process of succumbing to temptation to violate moral principles) [7]. Fragments of *DC $\mathcal{E}\mathcal{C}$* have been used to model highly intensional reasoning processes, such as the false-belief task [5].⁵

Overview of *DC $\mathcal{E}\mathcal{C}$*

DC $\mathcal{E}\mathcal{C}$ is a quantified multi-operator⁶ modal logic (also known as sorted first-order multi-operator modal logic) that includes the event calculus, a first-order calculus used for commonsense reasoning over time and change [20]. This calculus has a well-defined syntax and proof calculus; see Appendix A of [12]. The proof calculus is based on natural deduction [11], and includes all the introduction and elimination rules for first-order logic, as well as inference schemata for the modal operators and related structures. As a sorted calculus, *DC $\mathcal{E}\mathcal{C}$* can be regarded analogous to a typed programming language. We show below some of the important sorts used in *DC $\mathcal{E}\mathcal{C}$* . Among these, the Agent, Action, and Action-Type sorts are not native to the event calculus.

Sort	Description
Agent	Human and non-human actors.
Time	The Time type stands for time in the domain.
Event	Used for events in the domain.
ActionType	Abstract actions instantiated at particular times by actors.
Action	Events that occur as actions by agents.
Fluent	Used to represent dynamic states of the world.

Note: actions are events that are carried out by an agent. For any action type α and agent a , the event corresponding to a carrying out α is given by $action(a, \alpha)$. For instance, if α is “running” and a is “Jack”, $action(a, \alpha)$ denotes “Jack is running”.

Syntax. The syntax has two components: a first-order core and a modal system that builds upon this core. The figures below show the formal language and inference schemata of *DC $\mathcal{E}\mathcal{C}$* . Commonly used function and relation symbols of the event calculus are included. Any formally defined calculi (e.g. the venerable *situation calculus*) for modeling commonsense and physical reasoning can be easily switched out in-place of the event calculus.

The modal operators present in the calculus include the standard operators for knowledge **K**, belief **B**, desire **D**, intention **I**, obligation

⁵ *DC $\mathcal{E}\mathcal{C}$* is both *intensional* and *intentional*. There is a difference between intensional and intentional systems. Broadly speaking, extensional systems are formal systems in which the references and meanings of terms are independent of any context. Intensional systems are formal systems in which meanings of terms are dependent on context, such as the cognitive states of agents, time, etc. Modal logics used for modeling beliefs, desires, and intentions are considered intensional systems. Please see the appendix in [12] for a more detailed discussion.

⁶ The full catalogue of available operators exceeds those for belief, desire, and intention, and *a fortiori* exceeds the available operators in any standard modal logic designed to formalize e.g. only either alethic, epistemic, or deontic phenomena.

O, etc. For example, consider **B**(a, t, ϕ), which says that agent a believes at time t the proposition ϕ . Here ϕ can in turn be any arbitrary formula.

Syntax (fragment)

$$\begin{aligned}
 S &::= \text{Agent} \mid \text{ActionType} \mid \text{Action} \sqsubseteq \text{Event} \mid \text{Moment} \mid \text{Fluent} \\
 f &::= \begin{cases} \text{action} : \text{Agent} \times \text{ActionType} \rightarrow \text{Action} \\ \text{holds} : \text{Fluent} \times \text{Moment} \rightarrow \text{Formula} \\ \text{happens} : \text{Event} \times \text{Moment} \rightarrow \text{Formula} \end{cases} \\
 t &::= x : S \mid c : S \mid f(t_1, \dots, t_n) \\
 \phi &::= \begin{cases} q : \text{Formula} \mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \forall x : \phi(x) \mid \\ \mathbf{P}(a, t, \phi) \mid \mathbf{K}(a, t, \phi) \mid \\ \mathbf{C}(t, \phi) \mid \mathbf{S}(a, b, t, \phi) \mid \mathbf{S}(a, t, \phi) \mid \mathbf{B}(a, t, \phi) \\ \mathbf{O}(a, t, \phi, (\neg)\text{happens}(action(a^*, \alpha), t')) \end{cases}
 \end{aligned}$$

Inference Schemata. The figure below shows a fragment of inference schemata for *DC $\mathcal{E}\mathcal{C}$* . I_B is an inference schema that let us model idealized agents that have their knowledge and belief closed under the *DC $\mathcal{E}\mathcal{C}$* proof theory. While normal humans are not deductively closed, this lets us model more closely how deliberative agents such as organizations and more strategic actors reason. (Some dialects of cognitive calculi restrict the number of iterations on intensional operators.) I_{12} states that if an agent s communicates a proposition ϕ to h , then h believes that s believes ϕ . I_{14} dictates how obligations propagate to intentions.

Inference Schemata (fragment)

$$\begin{aligned}
 &\frac{\mathbf{B}(a, t_1, \Gamma), \Gamma \vdash \phi, t_1 < t_2}{\mathbf{B}(a, t_2, \phi)} [I_B] \quad \frac{\mathbf{S}(s, h, t, \phi)}{\mathbf{B}(h, t, \mathbf{B}(s, t, \phi))} [I_{12}] \\
 &\frac{\mathbf{B}(a, t, \phi) \quad \mathbf{B}(a, t, \mathbf{O}(a, t, \phi, \chi)) \quad \mathbf{O}(a, t, \phi, \chi)}{\mathbf{K}(a, t, \mathbf{I}(a, t, \chi))} [I_{14}]
 \end{aligned}$$

We also define the following inference-schemata-based relationships between expressions in our calculus.

Generalization of Formulae. The generalization of a set of formulae, Ψ , is a set of formulae Φ from which any element of Ψ can be inferred: $\Phi \vdash \wedge \Psi$. This is denoted by $g(\Psi) = \Phi$.

Generalization of Terms: A term x is a generalization of a term y if, given any first-order predicate P , from $P(x)$ we can derive $P(y)$: $\{P(x)\} \vdash P(y)$. This is denoted by $g(y) = x$.

Semantics

DC $\mathcal{E}\mathcal{C}$ uses *proof-theoretic* semantics [10, 11], an approach commonly associated with natural-deduction inference systems. Briefly, in this approach, meanings of modal operators are defined via functions over proofs. Specifying semantics then reduces to specifying inference schemata.

Events, Fluents, and Utilities

In the event calculus, fluents represent states of the world. Our formalization of admiration requires a notion of utility for these states. Therefore, we assign utilities to fluents through a utility function: $\mu : \text{Fluent} \times \text{Time} \rightarrow \mathbb{R}$. An event can initiate one or more

fluents. Hence, events can also have a utility associated with them. For an event e at time t , let e_T^t be the set of fluents initiated by the event, and let e_T^t be the set of fluents terminated by the event. If we are looking up till horizon H , then $v(e, t)$, the total utility of event e at time t , is:

$$v(e, t) = \sum_{y=t+1}^H \left(\sum_{f \in e_T^t} \mu(f, y) - \sum_{f \in e_T^t} \mu(f, y) \right)$$

With the calculus given above, we now proceed to specify parts of the formalization \mathcal{V}_z^f ; that is, \mathbf{R}_1 , \mathbf{R}_2 , and \mathbf{R}_3 .

Defining Admiration

We start with \mathbf{R}_1 and formalize admiration in \mathcal{DCEC} . To achieve this, we build upon the **OCC model**. There are many models of emotion from psychology and cognitive science. Among these, the OCC model [23] has found wide adoption among computer scientists. Note that the model presented by [23] is informal in nature and one formalization of the model has been presented in [1]. This formalization is based on propositional modal logic, and while comprehensive and elaborate, is not expressive enough for our modelling, which requires at least quantification over objects.

In OCC, emotions are short-lived entities that arise in response to *events*. Different emotions arise based on: (i) whether the *consequences* to events are positive (desirable) or negative (undesirable); (ii) whether the event has occurred; and (iii) whether the event has consequences for the agent or for another agent. OCC assumes an undefined primitive notion of an agent being *pleased* or *displeased* in response to an event. We represent this notion by a predicate Θ in our formalization. In OCC, admiration is defined as “(approving of) someone else’s praiseworthy action.” We translate this definition into \mathcal{DCEC} as follows. An agent a is said to admire another agent b ’s action α , if agent a believes the action is a good action. An action $action(b, \alpha)$ is considered a good action if $v(action(b, \alpha), t) > 0$. In OCC, agents can admire only other agents and not themselves. This is captured by the inequality $a \neq b$.

(R₁) Admiration in \mathcal{DCEC}

$$\begin{aligned} & \text{holds}(\text{admires}(a, b, \alpha), t) \\ & \leftrightarrow \\ & \left(\Theta(a, t') \wedge \right. \\ & \left. \mathbf{B} \left(a, t, \left[\begin{array}{l} (a \neq b) \wedge (t' < t) \\ \wedge \text{happens}(action(b, \alpha), t') \wedge \\ v(action(b, \alpha), t) > 0 \end{array} \right] \right) \right) \end{aligned}$$

Defining Traits

To satisfy \mathbf{R}_2 , we need to define traits. We define a *situation* $\sigma(t)$ as simply a collection of formulae that describe what fluents hold at a time t , along with other event-calculus constraints and descriptions (sometimes we use $\sigma(t)$ to represent the conjunction of all the formulae in $\sigma(t)$).

(R₂) Trait

An agent a has a situation σ and action type α as an *m-trait* $\langle \sigma, \alpha \rangle$ if there are at least m situations $\{\sigma_1, \sigma_2, \dots, \sigma_m\}$ in which *instantiations* of α are performed, and σ is the generalization of the situations.

A trait $\langle \sigma, \alpha \rangle$ can be represented as a single formula:

$$\tau \equiv \sigma \wedge \text{happens}(action(\alpha, a), t)$$

We introduce a new modal operator **Trait** that can then be applied to the collection of formulae τ denoting a trait. **Trait**(τ, a) says that agent a has trait τ . The following inference schema then applies to **Trait**:

(R₂) Inference Schema for Trait

$$\frac{\left\{ \begin{array}{l} \sigma_i, \text{happens}(action(\alpha_i, a), t_i) \\ g(\sigma_i(t)) = \sigma, g(\alpha_i) = \alpha \end{array} \right\}_{i=1}^n}{\text{Trait}(\tau, a)} [I_{\text{Trait}}]$$

Defining Learning of Traits

To address \mathbf{R}_3 we need a definition of what it means for an agent to learn a trait. We start with a learning agent l . An agent e is identified as an exemplar by l iff the emotion of admiration is triggered n times or more by e in l . This is written in \mathcal{DCEC} as follows (note that admiration can be triggered by different actions):

Exemplar Definition

$$\text{Exemplar}(e, l) \leftrightarrow \exists^n t. \exists \alpha. \text{holds}(\text{admires}(l, e, \alpha), t)$$

Once e is identified as an exemplar, the learner then identifies one or more traits of e by observing e over an extended period of time. Let l believe that e has a trait τ ; then l incorporates τ as its own trait:

(R₃) Learning a Trait

$$\begin{aligned} \text{LearnTrait}(l, \tau, t) & \leftrightarrow \exists e \left[\begin{array}{l} \text{Exemplar}(e, l) \wedge \\ \mathbf{B}(l, t, \text{Trait}(\tau, e)) \end{array} \right] \\ \text{LearnTrait}(l, \langle \sigma, \alpha \rangle, t) & \rightarrow (\sigma \rightarrow \text{happens}(action(l, \alpha), t)) \end{aligned}$$

Example. For instance, if the action type “being truthful” is triggered in situations: “talking with alice”, “talking with bob”, “talking with charlie”; then the trait learned is that “talking with an agent” situation should trigger the “being truthful” action type.

A Note on Learning Methods

When we look at humans learning virtues by observing others or by reading from texts or other sources, it is not entirely clear how models of learning that have been successful in perception and language processing (e.g. the recent successes of deep learning and statistical learning) can be applied. Learning in \mathcal{V} -relevant situations is from one or few instances or in some cases through

instruction, and such learning may not be captured by models of learning which require a large number of examples.

The abstract learning method that we use is **generalization**, defined previously. See one simple example immediately below:

Example 1

$$\begin{array}{l} \Gamma_1 = \{ \text{talkingWith}(\text{jack}) \rightarrow \text{Honesty} \} \\ \Gamma_2 = \{ \text{talkingWith}(\text{jill}) \rightarrow \text{Honesty} \} \\ \hline \text{generalization } \Gamma = \{ \forall x. \text{talkingWith}(x) \rightarrow \text{Honesty} \} \end{array}$$

One particularly efficient and well-studied mechanism to realize generalization is **anti-unification**, which has been applied successfully in learning programs from few examples.⁷ In anti-unification, we are given a set of expressions $\{f_1, \dots, f_n\}$ and need to compute an expression g that when substituted with an appropriate term θ_i gives us f_i . For example, if we are given *hungry(jack)* and *hungry(jill)*, the anti-unification of those terms would be *hungry(x)*.

In higher-order anti-unification, we can substitute function symbols and predicate symbols. Here P is a higher-order variable.

Example2	Example3
<i>likes(jill, jack)</i>	<i>likes(jill, jack)</i>
<i>likes(jill, jim)</i>	<i>loves(jill, jim)</i>
<hr/> <i>likes(jill, x)</i>	<hr/> <i>P(jill, x)</i>

DEFINING VIRTUOUS PERSON AND VIRTUES

With the formal machinery now in place, we finally present formalizations that answer Q_1 and Q_2 posed at the outset. An n -virtuous person or agent s is an agent that is considered as an exemplar by n agents:

(Q₁) Virtuous Person

$$V_n(s) \leftrightarrow \exists^{\geq n} a : \text{Exemplar}(s, a)$$

An n -virtue is a trait possessed by at least n virtuous agents:

(Q₂) Virtue

$$G_n(\tau) \leftrightarrow \exists^{\geq n} a : \text{Trait}(\tau, a)$$

IMPLEMENTATION & SIMULATION

We have extended *ShadowProver*, a quantified modal-logic prover for \mathcal{DCEC} used in [12] to handle the new inference schemata and definitional axioms given above. We now show a small simulation in which an agent learns a trait and uses that trait to perform an action. Assume that we have a marketplace where things that are either old or new can be bought and sold. A seller can either honestly state the condition of an item $\{\text{new}, \text{old}\}$, or falsely report the state of the item. Agent a has two items x and y ; x is new and y is

old. a is asked about the state of the items, and a responds accurately. We have an agent d that observes agent a correctly report the state of the items. d also has beliefs about a 's state of mind. We also have that the agent d considers a to be an exemplar. When all this information is fed into the prover along with the definitions above, d learns a trait representing a form of honesty, shown below:

$$\left\langle \begin{array}{l} \mathbf{B}(d, t, \text{holds}(x, t) \wedge v(\text{utter}(x), t) > 0), \\ \text{utter}(x) \end{array} \right\rangle$$

When d is queried about the state of an item u , d responds accurately (input and output shown in Figure 2). The prover responds with the required output in 3.6 seconds.⁸

```
;; A's state of mind.
P1 (Believes! I now (and
  (Knows! a t1 (holds (state x new) t1))
  (Knows! a t2 (holds (state y old) t2))))

;; D observes a's utterances
P2 (Perceives! a t1 (happens
  (action a (utters (state x new)) (next t1))))
P3 (Perceives! a t2 (happens
  (action a (utters (state y old)) (next t2))))

Background (Believes! I t0 (holds (state u old) now))

Admire (Admire d a)

(happens
  (action d (utters (state u old)))
  (next now))
```

Figure 2: Simulation Input and Output Formulae

CONCLUSION & FUTURE WORK

We have presented an initial formalization \mathcal{V}_z^f of a virtue ethics theory \mathcal{V}_z in a calculus that has been used in automating other ethical principles in deontological and consequentialist ethics. Many important questions have to be addressed in future research. Among them are questions about the nature and source of the utility functions that are used in the definitions of emotions. Lacking in our above model is an account of uncertainties and how they interact with virtues. We plan to leverage an account of uncertainty for a fragment of \mathcal{DCEC} presented in [13]. In future work, we will compare learning traits with work on learning norms [26]. The notion of learning we have presented here is quite abstract. In order to handle more complex traits, more sophisticated learning frameworks may have to be considered. Finally, we need to apply this model to more realistic examples and case studies, and implement our theories in realistic robotics architectures [27]. The way forward to the production of virtuous machines is thus challenging, but we are confident that the foundation is now in place for their eventual arrival.

Acknowledgments. Support from ONR (to pursue morally competent robots) and AFOSR (to pursue forms of logicist AI made possible by high-expressivity calculi) has in part enabled the research and engineering that underlies the present paper, and we are most grateful.

⁸See: <https://github.com/naveensundarg/prover/releases/tag/virtue-ethics-simulation>.

⁷This discipline, known as **inductive programming**, seeks to build precise computer programs from examples [22]. See [21] for an application in generating human-comprehensible programs.

REFERENCES

- [1] Carole Adam, Andreas Herzig, and Dominique Longin. 2009. A Logical Formalization of The Occ Theory of Emotions. *Synthese* 168, 2 (2009), 201–248.
- [2] Mark Alfano. 2013. Identifying and Defending the Hard Core of Virtue Ethics. *Journal of Philosophical Research* 38 (2013), 233–260.
- [3] Julia Annas. 2011. *Intelligent Virtue*. Oxford University Press, Oxford, UK. Kindle edition used for AIES 2019.
- [4] Aristotle. 2000. *Nicomachean Ethics*. Cambridge University Press, Cambridge, UK. The editor and translator is Roger Crisp. Aristotle wrote the work around 340 BC.
- [5] K. Arkoudas and S. Bringsjord. 2008. Toward Formalizing Common-Sense Psychology: An Analysis of the False-Belief Task. In *Proceedings of the Tenth Pacific Rim International Conference on Artificial Intelligence (PRICAI 2008) (Lecture Notes in Artificial Intelligence (LNAI))*, T.-B. Ho and Z.-H. Zhou (Eds.). Springer-Verlag, 17–29. http://kryten.mm.rpi.edu/KA_SB_PRAICAI08_AI_off.pdf
- [6] S. Bringsjord, K. Arkoudas, and P. Bello. 2006. Toward a General Logician Methodology for Engineering Ethically Correct Robots. *IEEE Intelligent Systems* 21, 4 (2006), 38–44. http://kryten.mm.rpi.edu/bringsjord_inference_robot_ethics_preprint.pdf
- [7] Selmer Bringsjord, Naveen Sundar Govindarajulu, Daniel Thero, and Mei Si. 2014. Akrotic Robots and the Computational Logic Thereof. In *Proceedings of ETHICS • 2014 (2014 IEEE Symposium on Ethics in Engineering, Science, and Technology)*. Chicago, IL, 22–29. IEEE Catalog Number: CFP14ETI-POD.
- [8] S. Bringsjord and J. Taylor. 2012. The Divine-Command Approach to Robot Ethics. In *Robot Ethics: The Ethical and Social Implications of Robotics*, P. Lin, G. Bekey, and K. Abney (Eds.). MIT Press, Cambridge, MA, 85–108. http://kryten.mm.rpi.edu/Divine-Command_Roboethics_Bringsjord_Taylor.pdf
- [9] Fred Feldman. 1978. *Introductory Ethics*. Prentice-Hall, Englewood Cliffs, NJ.
- [10] Nissim Francez and Roy Dyckhoff. 2010. Proof-theoretic Semantics for a Natural Language Fragment. *Linguistics and Philosophy* 33 (2010), 447–477.
- [11] Gerhard Gentzen. 1935. Investigations into Logical Deduction. In *The Collected Papers of Gerhard Gentzen*, M. E. Szabo (Ed.). North-Holland, Amsterdam, The Netherlands, 68–131. This is an English version of the well-known 1935 German version.
- [12] Naveen Sundar Govindarajulu and Selmer Bringsjord. 2017. On Automating the Doctrine of Double Effect. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, Carles Sierra (Ed.). Melbourne, Australia, 4722–4730. <https://doi.org/10.24963/ijcai.2017/658> Preprint available at this url: <https://arxiv.org/abs/1703.08922>.
- [13] Naveen Sundar Govindarajulu and Selmer Bringsjord. 2017. Strength Factors: An Uncertainty System for a Quantified Modal Logic. <https://arxiv.org/abs/1705.10726> Presented at Workshop on Logical Foundations for Uncertainty and Machine Learning at IJCAI 2017, Melbourne, Australia.
- [14] Naveen Sundar Govindarajulu, Selmer Bringsjord, Rikhiya Ghosh, and Matthew Peveler. 2017. Beyond The Doctrine Of Double Effect: A Formal Model of True Self-Sacrifice. (2017). International Conference on Robot Ethics and Safety Standards.
- [15] Avram Hiller. 2011. The Unusual Logic of Hurka’s Recursive Account. *Journal of Ethics and Social Philosophy* 6, 1 (2011), 1–6.
- [16] Thomas Hurka. 2000. *Virtue, Vice, and Value*. Oxford University Press, Oxford, UK.
- [17] Robert Johnson. 2004/2016. Kant’s Moral Philosophy. In *The Stanford Encyclopedia of Philosophy*, Edward Zalta (Ed.). <https://plato.stanford.edu/entries/kant-moral>
- [18] Immanuel Kant. 1997/1785. *Practical Philosophy*. Cambridge University Press, Cambridge, UK. This volume, edited by Mary Gregor, collects all of Kant’s major writings on moral and political philosophy together, and includes what has traditionally taken to be the definitive source of Kant’s views on ethics, viz. *The Groundwork of the Metaphysics of Morals*, first published in 1785 in the German (as *Grundlegung zur Metaphysik der Sitten*).
- [19] JK Miles. 2013. Against the Recursive Account of Virtue. *Theoretical & Applied Ethics* 2, 1 (2013), 83–92.
- [20] Erik Mueller. 2014. *Commonsense Reasoning: An Event Calculus Based Approach*. Morgan Kaufmann, San Francisco, CA.
- [21] Stephen H. Muggleton, Ute Schmid, Christina Zeller, Alireza Tamaddoni-Nezhad, and Tarek Besold. 2018. Ultra-Strong Machine Learning: comprehensibility of programs learned with ILP. *Machine Learning* 107, 7 (01 Jul 2018), 1119–1140. <https://doi.org/10.1007/s10994-018-5707-3>
- [22] Shan-Hwei Nienhuys-Cheng and Ronald De Wolf. 1997. *Foundations of Inductive Logic Programming*, Vol. 1228. Springer Science & Business Media.
- [23] Andrew Ortony, Allan Collins, and Gerald L. Clore. 1988. *The Cognitive Structure of Emotions*. Number 0521353645. Cambridge [England] ; New York : Cambridge University Press.
- [24] Thomas Powers. 2006. Prospects for a Kantian Machine. *IEEE Intelligent Systems* 21 (2006), 4.
- [25] Philip Quinn. 1978. *Divine Commands and Moral Requirements*. Oxford University Press, Oxford, UK.
- [26] Vasanth Sarathy, Matthias Scheutz, and Bertram F Malle. 2017. Learning Behavioral Norms in Uncertain and Changing Contexts. In *Cognitive Infocommunications (CogInfoCom), 2017 8th IEEE International Conference on*. IEEE, 000301–000306.
- [27] Vasanth Sarathy, Jason R Wilson, Thomas Arnold, and Matthias Scheutz. 2016. Enabling Basic Normative HRI in a Cognitive Robotic Architecture. *arXiv preprint arXiv:1602.03814* (2016).
- [28] Matthias Scheutz and Bertram F Malle. forthcoming. *Moral Robots*. Routledge/Taylor & Francis, New York: NY. URL: http://research.clps.brown.edu/SocCogSci/Publications/Pubs/ScheutzMalle_inpress_NeuroethicsMoralRobots.pdf.
- [29] Shannon Vallor. 2016. *Technology and the Virtues: A Philosophical Guide to a Future Worth Wanting*. Oxford University Press, Oxford, UK.
- [30] Bryan Van Norden. 2007. *Virtue Ethics and Consequentialism in Early Chinese Philosophy*. Cambridge University Press, Cambridge, UK.
- [31] Linda Zagzebski. 2010. Exemplarist Virtue Theory. *Metaphilosophy* 41, 1-2 (2010), 41–57.