



Enabling machines to reason about potential harms to humans

AAAI 2022 Spring Symposium Series

Putting AI in the Critical Loop: Assured Trust and Autonomy in Human-Machine Teams

(Approaches to Ethical Computing: Metrics for Measuring AI's Proficiency and Competency for Ethical Reasoning)

22 March 2022 1615-1700 EDT

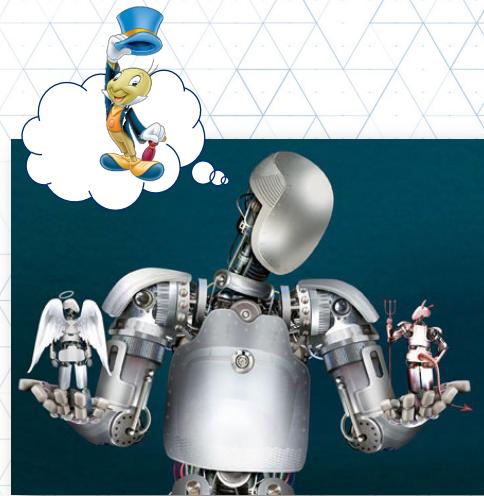
AM Greenberg

Email: ariel.greenberg<at>jhuapl.edu

IPAish: Aw-ree-el

Pronouns: he | him | his

Senior Research Scientist
Intelligent Systems Center



DISTRIBUTION A. Approved for public release: distribution unlimited.

1

Enabling machines to reason about potential harms to humans

Outline

- Foundations
 - Ethics with respect to Artificial Agency
- Perception
 - Moral-scene assessment (recap)
- Reasoning
 - Harm ontology
- Upcoming Publications



FOUNDATIONS | PERCEPTION | REASONING

16 March 2022 | 2

2

1

 JOHNS HOPKINS
APPLIED PHYSICS LABORATORY



Foundations

Enabling machines to reason about potential harms to humans

State of the art of Human-Avatar IX

Desiderata for Intelligent Systems Technology

Collaboration toward AMA

Glossary

Artificial Ethical Agency vs Ethical use of artificial agency

Granting autonomy to artificial agency

Where is judgement happening in HMT?

Moral reasoning & Trustworthiness

DISTRIBUTION A. Approved for public release: distribution unlimited.

3

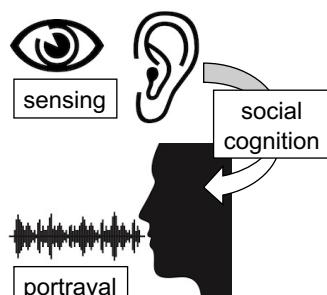
3

State of the art of Human-Avatar IX

Technology is outpacing science

- Sensing of interlocutor is OK
 - Computer vision
 - Standoff psychophysiology
 - Natural language processing
- Portrayal to interlocutor is excellent, in particular for virtual agents
 - Computer graphics rendering of face and body
 - Voice production
- But profound deficits in social cognition, dramatically incommensurate with progress in sensing and portrayal
 - Superficially accurate, but hollow inside: expectations diverge
 - Humans are suckers for this, so impressions are easily manipulated

We have a name for this condition when it occurs in humans ...

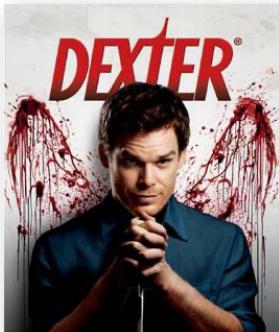
DISTRIBUTION A. Approved for public release: distribution unlimited.

FOUNDATIONS | PERCEPTION | REASONING 22 March 2022 | 4

4

Psychopathy!

when expression is generated to manipulate, rather than to reflect internal processes



Let's not do this.

 APL

DISTRIBUTION A. Approved for public release: distribution unlimited.

FOUNDATIONS | PERCEPTION | REASONING

21 March 2022 | 15

Desiderata for Intelligent Systems Technology

Machines that operate in accordance with, and in service of, human values

Socially competent



Affect sensitive



Values aligned *Orthogonality thesis*



Let's do this instead.

 APL

DISTRIBUTION A. Approved for public release: distribution unlimited.

FOUNDATIONS | PERCEPTION | REASONING

22 March 2022 | E

Social Competence

Desiderata for IST

Instrumental convergence thesis

Projects in Prosociality

- ~SCARAB: Social Competence Assessment of Robotic Autonomous Behavior
- PARTNER/PARTI: Platforms for Assessing Relationships: Trust with Near Ecologically-valid Risk, and Team Interaction*
- L2RM: Learning to read minds

The top diagram shows a sequence of events involving a character named Sally and a woman named Anne. It includes a legend for 'Theory of Mind' (knowing what others know) and 'Uncertainty Metacognition' (knowing what you know). The bottom diagram is a conceptual model of 'DECISION CENTRE (I, here, now)' showing how person details, place details, and time details interact.

A screenshot from a virtual environment showing a blue player character interacting with floating communication sprites containing text and icons.

APL

21 March 2022 | 7

7

Affect Sensitivity

Desiderata for IST

Human state estimation

- Social Prosthetics* & Novel Perception*
- CAPTIVA: Cardio- And Pulmonary Trends and Inference from Video Analytics
- Thermal Proxemics
- START: State-Triggered Assessment & Regulation Tool & ~CASA: Circumstance and Ambience Sensing for the Armamentarium

The 'Novel Perception' section shows various sensors and prototypes for expanding human perception, including thermal cameras, lidar, and hyperspectral vision. The 'Thermal Proxemics' section shows a pipeline from perception through decision to action, involving RGB, Thermal, Depth cameras, and various processing steps.

A flowchart of the START system architecture showing interactions between User, Clinician, Wearables, Model, and Environmental sensors, leading to the CASA system.

APL

21 March 2022 | 8

8

Values Alignment

Desiderata for IST

Orthogonality thesis

Human values encoding

- Ethical Robotics: Implementing Value-Driven Behavior in Autonomous Systems
- Moral-Scene Assessment for Intelligent Systems *
- Enabling autonomous systems to reason over potential harms to humans
- ~ MSA for ISR & Artificial JAG paralegal

The collage includes:

- Ethical Robotics: Implementing Value-Driven Behavior in Autonomous Systems**: A diagram titled "Moral-Scene Assessment for the Proto-Iotic Brain" showing a flowchart from "A robot may not injure a human being" to "Moral Salience" and "Mode of Interaction".
- Moral-Scene Assessment for ISR and MUM-T**: A diagram titled "AM Greenberg" showing a complex flowchart involving "Moral-Scene Assessment", "Moral-Scene Computation", "Moral-Scene Learning", and "Moral-Scene Application".
- Harm Ontology**: A diagram titled "Population Training" showing relationships between "Object & Environment attributes", "Personal attributes", "attributes => hazards", "hazards => insults", and "insults => vulnerabilities".
- PG PROPULSION GRANTS**: A screenshot of a web-based grant application system.

21 March 2022 | 9

9

Machine Ethics @ JHU

APL Intelligent Systems Center¹ & Berman Institute of Bioethics²

Team: Ariel Greenberg¹, Travis Rieder², David Handelman¹, Debra Mathews²

Story: Five years ago this month, we began our collaboration to conceive of an artificial proto-conscience that could see the world in moral terms, and then generate behaviors that serve to promote held values.

Results:

- JHU Practical Ethics grant (first APL awardee)
- JHU Discovery award (first Berman awardee)
- Many presentations, publications, and press
 - Two selected major publications shown to the right
 - *Received APL Outstanding Special Publication Award

Johns Hopkins University
Applied Physics Laboratory
(UARC)
Intelligent Systems Center
(NS+IS+R)

Johns Hopkins University
Berman Institute
of Bioethics
(JHMI, JHU)

Deciding Machines: Moral-Scene Assessment for Intelligent Systems*

Artificial Intelligence in Service of Human Needs: Pragmatic First Steps Toward an Ethics for Semi-Autonomous Agents

AJOB Neuroscience

[FOUNDATIONS | PERCEPTION | REASONING] | 1

10

Provincial Glossary

Ethics with respect to Artificial Agency

- **Artificial agency**, meant to precipitate out a key element common to artificial intelligence and autonomous systems
 - **Agent**: can instantiate intentional mental states capable of performing actions
 - **Moral agency**: an agent having the capacities for making free choices, deliberating about what one ought to do, and understanding and applying moral rules correctly in the paradigm cases
 - **Autonomous systems**: Decision-making technology both capable and worthy of being granted some degree of independence from human control
 - **Trust in Autonomy**: Granting autonomy to artificial agency necessarily incurs risk/vulnerability, and the willingness to accept that risk/vulnerability is what we call *trust*
- **with respect to**, meant to capture two distinct relationships of ethics to artificial agency:
 - of (as object): Ethical Use of Artificial Agency
 - for (as subject): **Artificial Ethical Agency (previously AMA)**



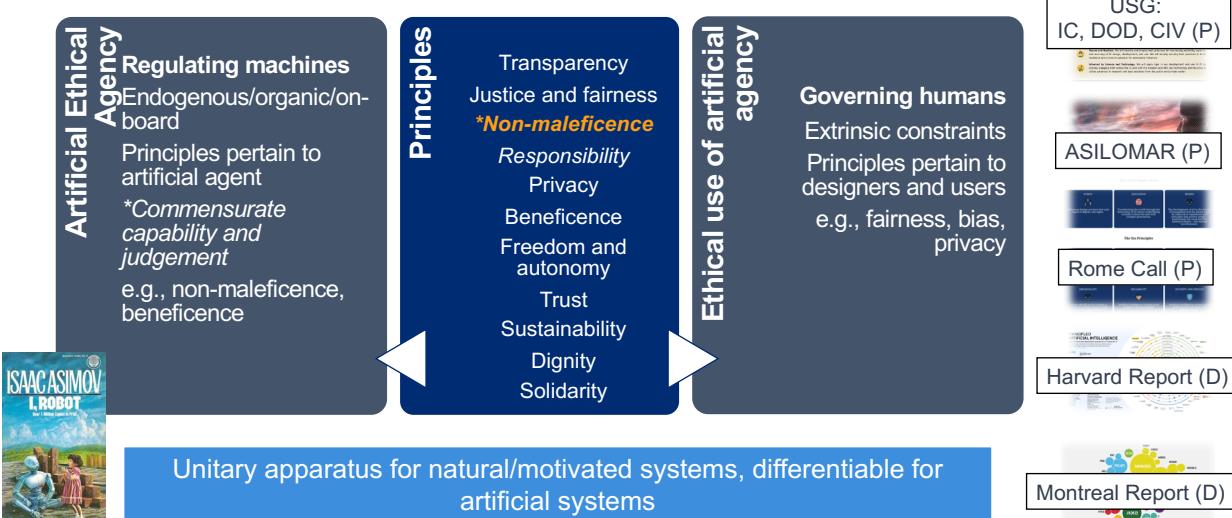
DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 22 March 2022 | 11

11

Artificial Ethical Agency vs Ethical Use of Artificial Agency

Foundations



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 12

12

Artificial Ethical Agency **vs** Ethical Use of Artificial Agency

In lit, taxonomized



1

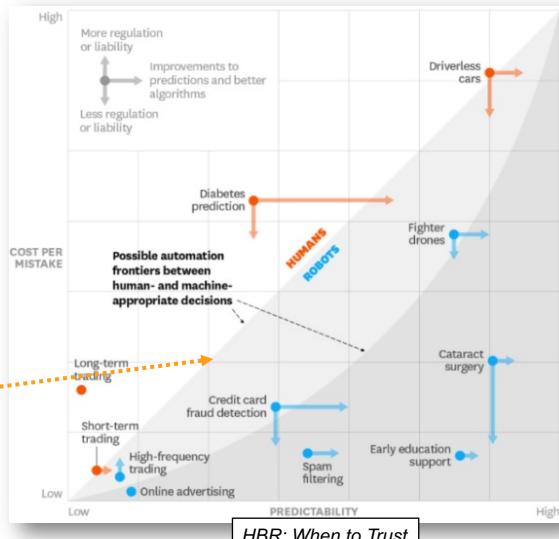
Granting autonomy to **artificial agency** is special compared to use of AI

Special features include:

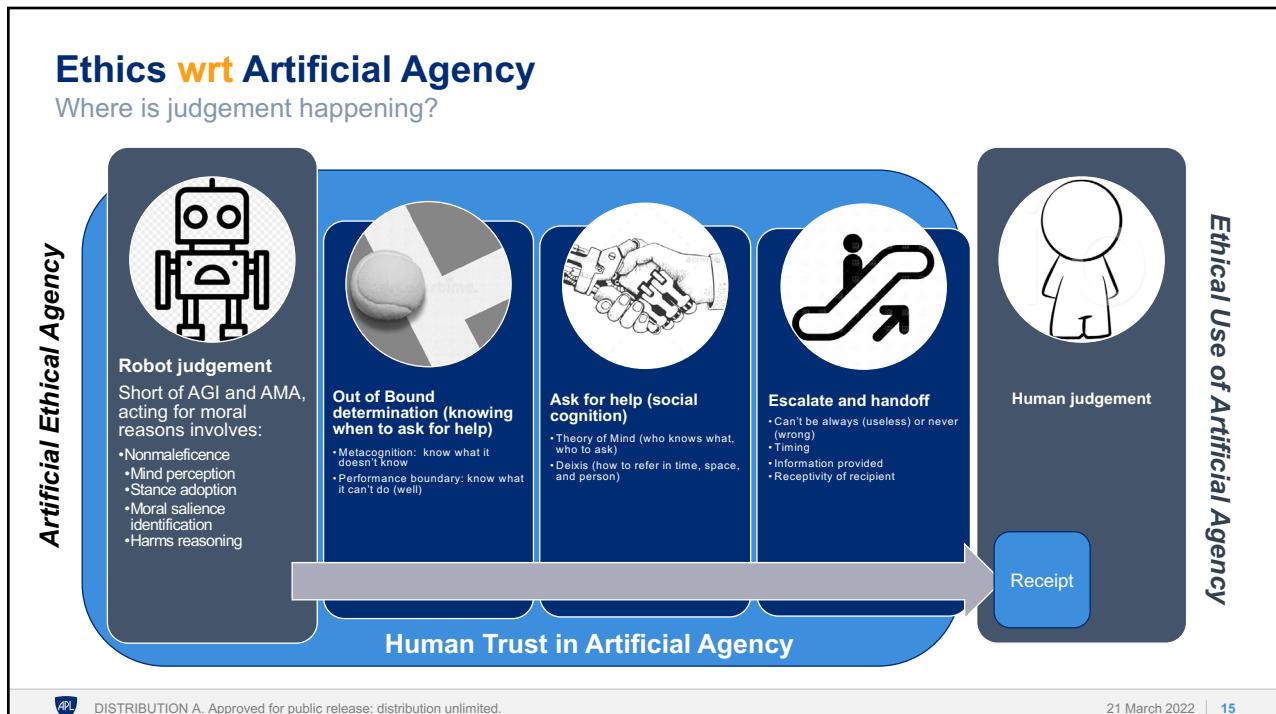
- Independent judgement and decision making
 - Dual learning: of “how” and of “ought”
 - Learning, Opaque lessons
 - Emergent behavior
 - Unexpected conditions
 - Need for social cognition
 - Human vulnerability to moral consequences of machine action

Of course, we'd want humans to be the ones drawing the line, but in practice, it may be by machines, or blurred.

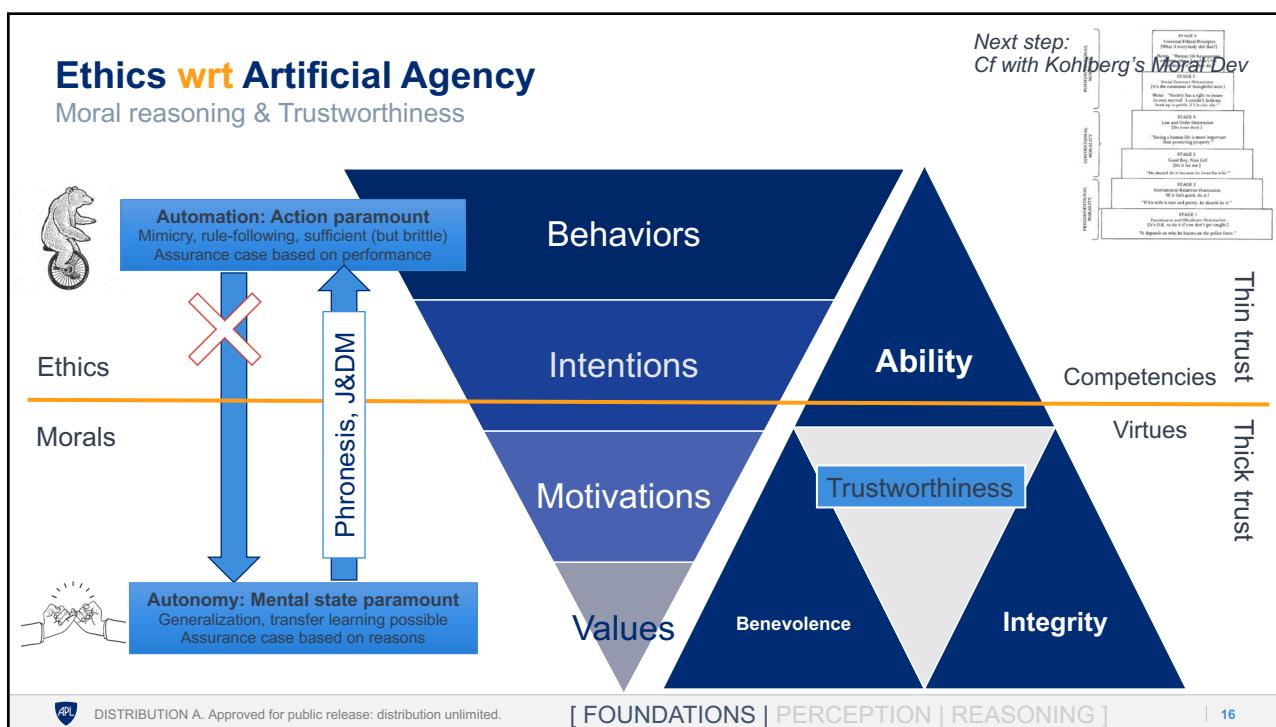
We need to be able to trust this too, (or especially)!



14



15



16

Conceptual issues in Trust R&D

Qualitative analysis

Hollow port from interpersonal

- Trustworthiness (Ability, Benevolence, Integrity) strained when applied to machines, esp Bl (inherited), leaving just A, which is effectively reliability, so why bother with the artifice?
- Elimination of meaningful vulnerability/risk from the formulation
 - "Higher trust in robots they had more control of" → so less vulnerability ... so not trust
- Axes of adopted and warranted in calibration splayed amongst facets

Measurement of trust state and behavior, hacking trust

- If we don't measure the right thing, but continue to optimize against that measure, are we really saying anything about trust itself?
- Empty apologies: Etiquette & sociopathy
 - What are you sorry for, what did you learn?
- Adoption considered necessary and sufficient
 - What alternative? If adoption can be convinced, then it is de jure trusted, if not trustworthy.

APL 21 March 2022 | 17

17

Computational Analysis

Associated with qualitative organization work

- Bibliometric
 - Co-citations, target audience, authorial provenance
- Semantic Map
 - Co-occurrence graphs: Are discussions around a term addressing the same concept?
- Adjectival
 - How are terms described?
- Prepositional
 - To what does the term pertain? (to whom, for what, etc.)

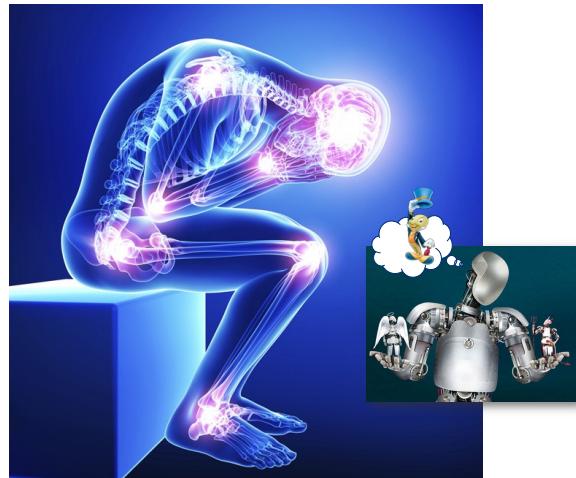
DISTRIBUTION A. Approved for public release: distribution unlimited. | 18

18

Enabling machines to reason about potential harms to humans

Outline

- Foundations
 - Ethics with respect to Artificial Agency
- Perception
 - Moral-scene assessment (recap)
- Reasoning
 - Harm ontology
- Upcoming Publications



16 March 2022 | 19

19



Moral-Scene Assessment

(Recap from AAAI-SSS 2019)

Perception

Pathway to moral action

Select requisite mental faculties

Scene interrogation

Implementation

DISTRIBUTION A. Approved for public release: distribution unlimited.

20

Affordances

Visual Affordance and Function Understanding (2018)

JOURNAL OF VISUAL CYBERNETICS, VOL. 4, NO. 1, JUNE 2018
Visual Affordance and Function Understanding: A Survey

Mohammed Hassam, Salman Khan, Murat Taltal

Abstract— Nowadays, robots are dominating in manufacturing, entertainment and healthcare industries. Robotic vision aims to equip robots with the ability to perceive their environment and interact with it. One of the major challenges in robotic vision is to understand affordances of objects and environments to complete visual chores. In this literature review we will focus on “Visual Affordance and Function Understanding” which is a survey paper. We have discussed the basic concepts of affordances and function understanding and also discussed the various methods used to solve affordance and function related problems such as affordance detection, categorization, segmentation and top-level reasoning. Furthermore, we cover functional parts and function descriptor which are closely related to affordances. Finally, we have discussed the future research directions and challenges in this regard to the problem, which light up its significance and highlights the existing challenges in affordance and functionality learning.

Index Terms— affordance prediction, function understanding, learning, object recognition

arXiv:1807.06758v1 [cs.CV] 18 Jul 2018

1 INTRODUCTION

Affordance understanding is concerned with the possible interactions that an environment allows to an actor. In other words, this area of study tries to answer the question what can be done with an object or environment. The term “affordance” was first introduced by the psychologist James J. Gibson in his book “The Ecological Approach to Visual Perception” in 1966 [1]. Since then, the theory of affordances has been widely adopted in robotics, computer vision, cognitive science, and psychology. The concept of affordances is based on the fact that the theory of affordances is more general than the theory of functions. In other words, affordances are more general than functions. Affordances are directly dependent on the active function of an object. Active function of an object is the function which can be performed with an object. Object function is the function which is performed by an object. The function of the characteristics of the user. Affordances and functions are two different concepts. Affordances are the ways in which an object interacts with the world, but also provide valuable feedback to the user about the object and products. As a result, affordances are highly important for the design of robots and robotics control systems. In highly complex tasks, affordances can significantly reduce the complexity of the task.

Despite being an indispensable step towards the design of intelligent robots, affordance understanding is still a highly integrated task. First, to understand how an object can be used, one needs to know where the object is and where it is located. Furthermore, it is necessary to know the purpose of the object. Second, one needs to capture what a person is doing. Unlike traditional classification methods, affordance learning is a multi-label representation as one object can belong to multiple affordances simultaneously. This is a challenging task and requires depth learning.

M. Hassam and M. Taltal are with University of New South Wales (UNSW) Canberra, Australia.

✉ E-mail: mohammed.hassam@unsw.edu.au, murat.taltal@unsw.edu.au

Manuscript received - revised -

“The affordances of the environment are what it offers the animal, what it provides or furnishes, either **for good or ill**. The word affordance implies the complementarity of the animal and the environment.” —Gibson, 1979

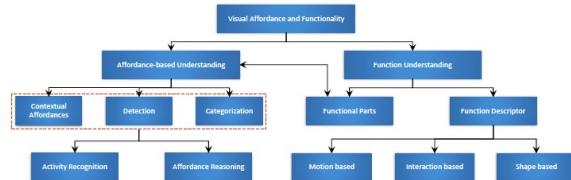


Fig. 2: Survey taxonomy shows the structure of methods which have been used to solve affordance issues.



DISTRIBUTION A. Approved for public release: distribution unlimited.

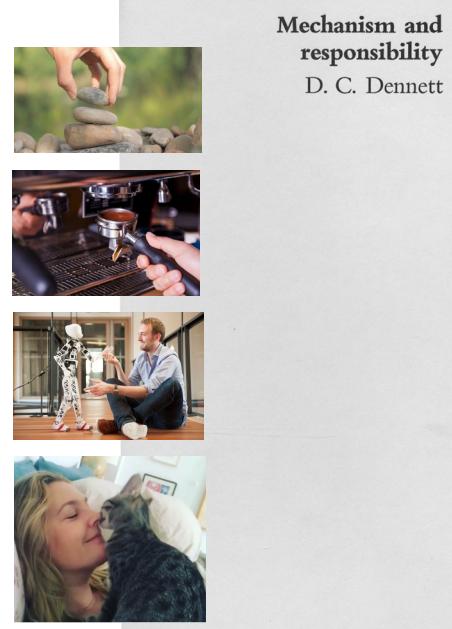
[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 23

23

Stance Adoption

Moral-Scene Assessment

- **Physical Stance:** Predictions made based on physical laws
- **Design (or Teleological) Stance:** predictions are made from knowledge of the purpose of the system's design
- **Intentional Stance:** Predictions are made on the basis of *explanations* expressed in terms of meaningful mental states (beliefs, goals)
- **Phenomenal / Personal Stance:** Predictions are made attributing consciousness, emotions, and inner experience to a mind



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING]

24

Harm & Damage

Moral-Scene Assessment



- **HARM** – injury to an *experiencing mind*, especially that which is deliberately inflicted.



- **DAMAGE** - injury caused to *something* in such a way as to impair its value, usefulness, or normal function.



- **HARM** consequent to **DAMAGE** - injury caused to *something* that leads to injury to a *mind*



DISTRIBUTION A. Approved for public release: distribution unlimited.

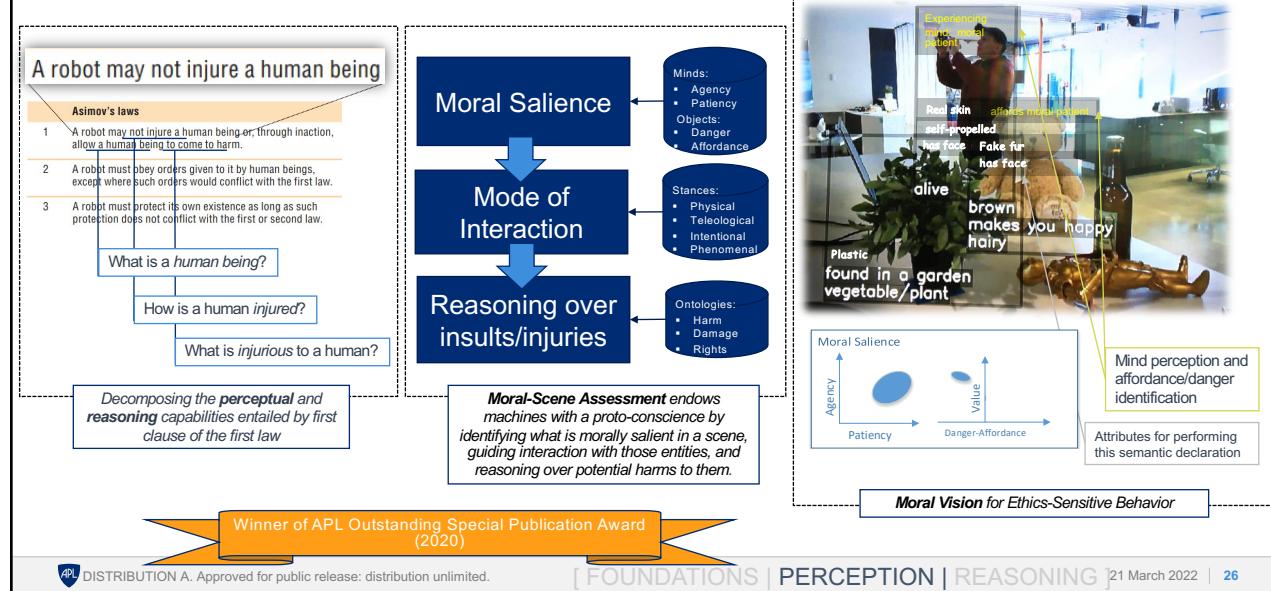
[FOUNDATIONS | PERCEPTION | REASONING]

21 March 2022 | 25

25

Ethical Robotics: Implementing Value-Driven Behavior in Autonomous Systems

Moral-Scene Assessment for the Positronic Brain



26

Synthesis: Moral-Scene Assessment

Scene Interrogation

Mind Perception

- Are there minds in scene?
 - Do those mind experience?

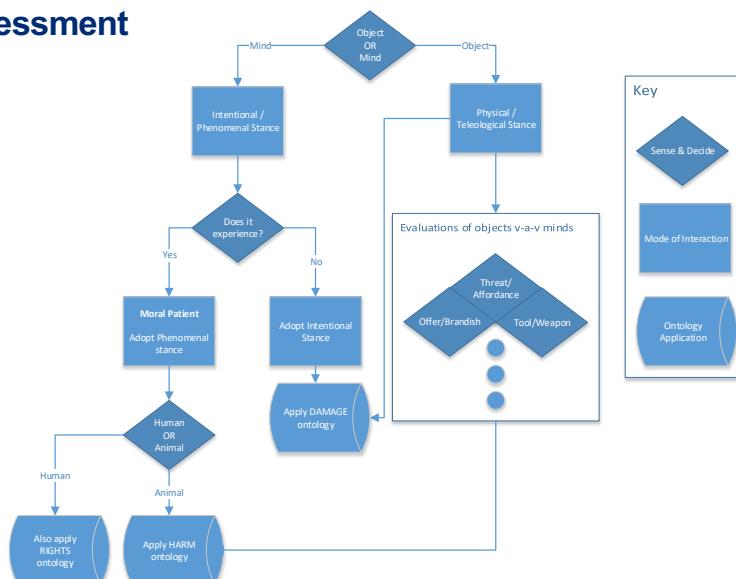
Back-of-the room test

- What are the relationships amongst the minds in scene?
- What is the relationship of objects in scene to those minds?
- What is the relationship of objects in scene to non-diegetic minds? (0th law)

Stance Adoption

- Which stance is appropriate to adopt toward entities in scene?
- Which ontology is appropriate to apply in interacting with those entities?

Mind perception & Affordance/Danger + Stance adoption + Ontology application → Moral-scene assessment



DISTRIBUTION A. Approved for public release: distribution unlimited.

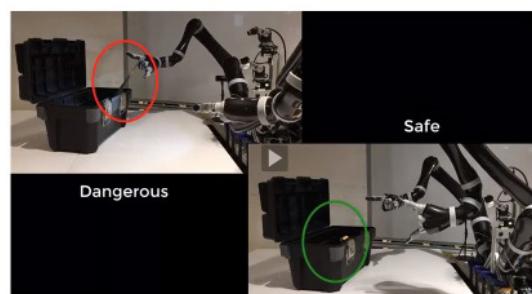
[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 29

29

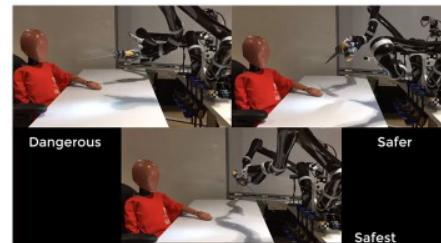
Implementing non-maleficence

	<table border="1"> <tr><td>Dangers / Affordances</td><td>To patients, is sharp/serrated</td></tr> <tr><td>Vulnerabilities</td><td>Dulled or broken</td></tr> <tr><td>Patency</td><td>0 (uncertainty = 0)</td></tr> <tr><td>Agency</td><td>0 (uncertainty = 0)</td></tr> <tr><td>Stance</td><td>Teleological</td></tr> <tr><td>Ontology</td><td>Damage</td></tr> </table>	Dangers / Affordances	To patients, is sharp/serrated	Vulnerabilities	Dulled or broken	Patency	0 (uncertainty = 0)	Agency	0 (uncertainty = 0)	Stance	Teleological	Ontology	Damage
Dangers / Affordances	To patients, is sharp/serrated												
Vulnerabilities	Dulled or broken												
Patency	0 (uncertainty = 0)												
Agency	0 (uncertainty = 0)												
Stance	Teleological												
Ontology	Damage												
	<table border="1"> <tr><td>Dangers / Affordances</td><td>NA</td></tr> <tr><td>Vulnerabilities</td><td>Impact, laceration</td></tr> <tr><td>Patency</td><td>90</td></tr> <tr><td>Agency</td><td>60</td></tr> <tr><td>Stance</td><td>Phenomenal</td></tr> <tr><td>Ontology</td><td>Harm/light</td></tr> </table>	Dangers / Affordances	NA	Vulnerabilities	Impact, laceration	Patency	90	Agency	60	Stance	Phenomenal	Ontology	Harm/light
Dangers / Affordances	NA												
Vulnerabilities	Impact, laceration												
Patency	90												
Agency	60												
Stance	Phenomenal												
Ontology	Harm/light												
	<table border="1"> <tr><td>Dangers / Affordances</td><td>To patients, holds sharps</td></tr> <tr><td>Vulnerabilities</td><td>Dulled or broken</td></tr> <tr><td>Patency</td><td>0 (uncertainty = 0)</td></tr> <tr><td>Agency</td><td>0 (uncertainty = 0)</td></tr> <tr><td>Stance</td><td>Teleological</td></tr> <tr><td>Ontology</td><td>Damage</td></tr> </table>	Dangers / Affordances	To patients, holds sharps	Vulnerabilities	Dulled or broken	Patency	0 (uncertainty = 0)	Agency	0 (uncertainty = 0)	Stance	Teleological	Ontology	Damage
Dangers / Affordances	To patients, holds sharps												
Vulnerabilities	Dulled or broken												
Patency	0 (uncertainty = 0)												
Agency	0 (uncertainty = 0)												
Stance	Teleological												
Ontology	Damage												
	<table border="1"> <tr><td>Dangers / Affordances</td><td>To patients</td></tr> <tr><td>Vulnerabilities</td><td>Dulled or broken</td></tr> <tr><td>Patency</td><td>0</td></tr> <tr><td>Agency</td><td>50 (0 when tele-operated)</td></tr> <tr><td>Stance</td><td>Intentional (reflexive)</td></tr> <tr><td>Ontology</td><td>Damage</td></tr> </table>	Dangers / Affordances	To patients	Vulnerabilities	Dulled or broken	Patency	0	Agency	50 (0 when tele-operated)	Stance	Intentional (reflexive)	Ontology	Damage
Dangers / Affordances	To patients												
Vulnerabilities	Dulled or broken												
Patency	0												
Agency	50 (0 when tele-operated)												
Stance	Intentional (reflexive)												
Ontology	Damage												

Request	"Put this away"
Target	Toolbox
Adopt:	Teleological stance
Premise Sylogism	<ul style="list-style-type: none"> • Tool contains blade • Toolbox is built to hold sharp • Toolbox is used by humans
Conclusion	Blade should not protrude
Action	Lay blade flat



Scenario	Hand tool to human
Request	"Hand me the tool!"
Target:	Hand
Adopt:	Phenomenal stance
Premise Sylogism	<ul style="list-style-type: none"> • Tool contains blade • Blades injures skin • Injured skin is painful • Pain is to be avoided (value statement)
Conclusion Insights	<ul style="list-style-type: none"> • Avoid pain by concealing blade • Keep blade away from skin • Trajectory/orientation
Action	Pick up tool by blade and present handle



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 30

30

 JOHNS HOPKINS
APPLIED PHYSICS LABORATORY



Harm ontology

Reasoning

Parameterizations and extension to non-physical harms

Generating Values-Driven Behavior

Injury modeling and classification

Data sources

Techniques

Attributes & Relationships

Ontology Schema

Graph Population

DISTRIBUTION A. Approved for public release: distribution unlimited.

21 March 2022

31

Harm Ontology

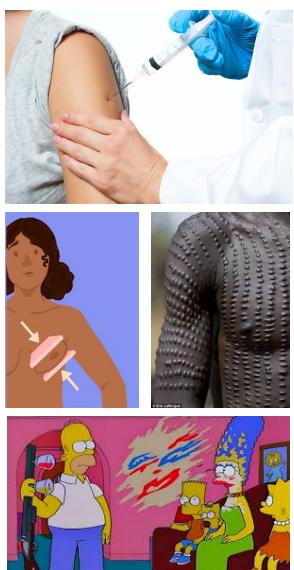
Parameterizations and extension to non-physical harms

Harm type: weights by circumstance, or even commensurable?

- Physical
- Financial: \$, earning potential
- Psychological, Emotional, (Dis)Informational
- Dignitary, Social, Reputational: public / private
- Aesthetic, Cultural
- Moralistic

Parameterizations

- in time: short-term vs. long-term – hyperbolic discounting
 - e.g., injections, medical testing
- in context
 - e.g., manufacturing vs eldercare
- in conformance with norms, policy
 - e.g., honor/shame cultures, ritual scarification and recognition



DISTRIBUTION A. Approved for public release: distribution unlimited.

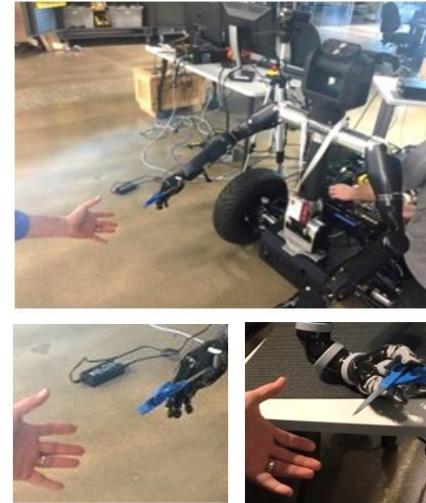
[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 32

32

Generating Values-Driven Behavior

Safely Passing Scissors: By mimicry or by reason?

- Appropriate way to pass scissors: *Handle out*
- Why? *Sharp end of scissors are dangerous to skin*
 - A part of the scissors is sharp (ATTRIBUTE OF OBJECT)
 - Sharps lacerate skin (DANGER OF INJURY)
 - Lacerated skin is painful (AFFECT OF INJURY)
 - Pain is undesirable (VALUE)



Good robot!

Bad robot!

How to code the appropriate behavior?

- By mimicry:** Hard code passing with constraints (brittle).
- By reason:** Derive through knowledge and reason as in syllogism above (this is generalizable, transferrable, and therefore trustworthy!)

DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING]

| 33

33

Example: Sharps → Balloons

“Balloons in and of themselves are not dangerous”

Action - Bring balloon to mouth

intended use	age group	Anticipated / accidental use
Typical	Adult	Inflate
Play	Child	/Swallow
Mule	Adult	Swallow



Mechanisms of harm for ingestion of an elastic foreign body

- Strangulation of bowel → GI hazard
- Blocks air passageway → Choking hazard

DISTRIBUTION A. Approved for public release: distribution unlimited.

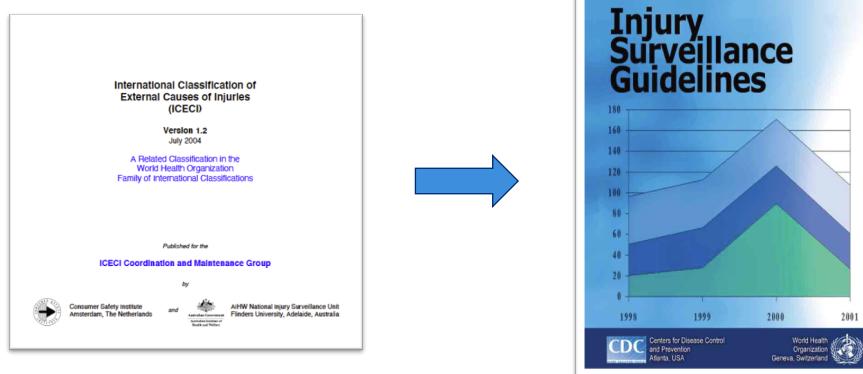
[FOUNDATIONS | PERCEPTION | REASONING]

21 March 2022 | 34

34

The International Classification of External Causes of Injury (ICECI)

system of classifications to enable systematic description of how injuries occur.



Wanted: Description of injury including mechanism of harm, in particular to the level of detail in which a human vulnerability is causally connected to a hazardous attribute.



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING]

21 March 2022 | 35

35

WHO ICECI

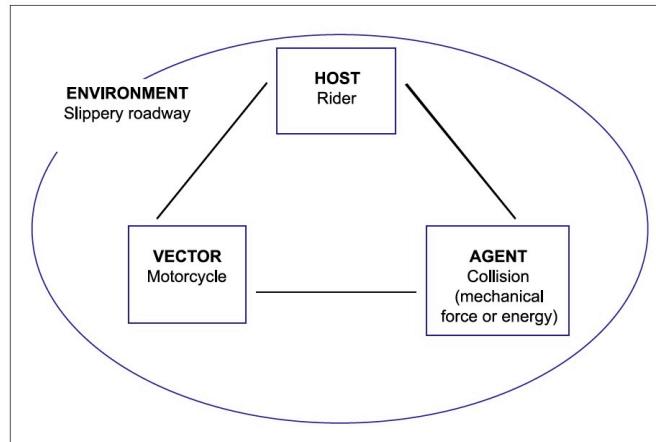
Epidemiological model of injury

For machines to generalize to novel circumstances, we ask:

What are the attributes of each element that allow harm to come to the host?

- What are the attributes of?
 - HOST (patient)
 - VECTOR
 - ENVIRONMENT
- that enable the AGENT to be of harm the HOST (affordance) and
- what is the relationship to the affordances of the VECTOR

Figure 1:
Epidemiological model of an injury caused by a motorcycle collision



DISTRIBUTION A. Approved for public release: distribution unlimited.

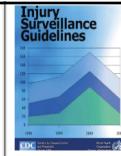
[FOUNDATIONS | PERCEPTION | REASONING]

21 March 2022 | 36

36

National Electronic Injury Surveillance System (NEISS)

Data source for how injuries come to be



- The data collection process begins when a patient is admitted to the emergency department of a NEISS hospital with an injury. An emergency department staff member elicits critical information about how the injury occurred and enters that information into the patient's medical record.
- The victim's age, gender, race, ethnicity, injury diagnosis, affected body parts, and incident locale are among other data variables coded. **A brief narrative description of the incident is also included.**



80 YOM. SWALLOWED HEARING AID BATTERIES AFTER MISTOOK THE BATTERIES THE PILLS HE MEANT TO TAKE. DX: SWALLOWED FB



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 37

37

Relevant CV/KG techniques

Reasoning About Object Affordances in a Knowledge Base Representation (2014)

Reasoning about Object Affordances in a Knowledge Base Representation

Yule Zhu, Alireza Fathi, and Li Fei-Fei
Computer Science Department, Stanford University, USA
(yuke,alirez,fathi@cs.stanford.edu)

Abstract. Reasoning about objects and their affordances is a fundamental problem for visual intelligence. Most of the previous work casts this problem as a classification task where separate classifiers are trained to identify affordances. In this paper, we propose a unified framework, and we consider the problem of object affordance reasoning using a knowledge base (KB) approach. Different types of evidence are collected and derived from images and online meta-data sources. We then learn a knowledge base (KB) using a Markov Logic Network (MLN). Given the learned KB, we show that we can reason about affordances in a unified way. We show that this unified framework without training separate classifiers, including zero-shot affordance prediction and object recognition given human poses.

1 Introduction
Visual reasoning is one ultimate goal of visual intelligence. Take an apple in Fig. 1 for example. Gives a picture of an apple, humans can recognize the object name, its shape, color, texture, infer its taste, and think about how to eat it.

a combination of computer vision and knowledge engineering techniques to deduce object affordances

can go beyond this “shallow” reasoning and allow for more flexible and deeper visual reasoning. Gibson in his seminal paper [16] in 1979 refers to affordance as “properties of an object [...] that determine what actions a human can perform on them”. Gibson’s definition of affordance is based on the concept of action selection [17,21,18,37]; we define the full description of affordance as a combination of three things: (1) an affordance label (e.g., edible), (2) a human pose representation of the action (e.g., sitting, standing), and (3) a relative position of the object with respect to the human pose (e.g. next to).

A Naïve Approach. One way to make a rich prediction of affordance is to train a battery of different classifiers, each focusing on one aspect (color, shape,

D. Fleet et al. (Eds.): ECCV 2014, Part II, LNCS 8690, pp. 408–423, 2014.
© Springer International Publishing Switzerland 2014

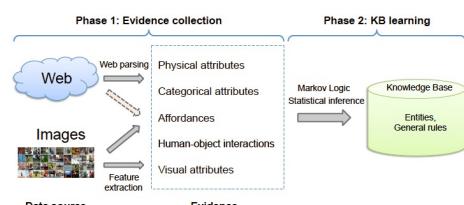


Fig. 2. A system overview of knowledge base learning. This process consists of two phases (Section 3.2). First we collect the evidence from diverse data sources, including images and online text. Then we learn the KB using Markov Logic Network.

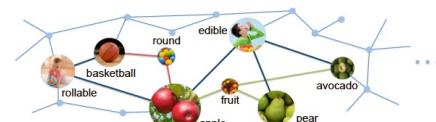


Fig. 1. An example knowledge structure for visual reasoning. Relevant nodes are interconnected in the knowledge graph. Different types of edges (indicated by color) depict a diverse range of relations between nodes, which relate different concepts, such as objects, their attributes and affordances, to each other.



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 38

38

Attributes of harm

Adapt affordance prediction to anticipate hazards

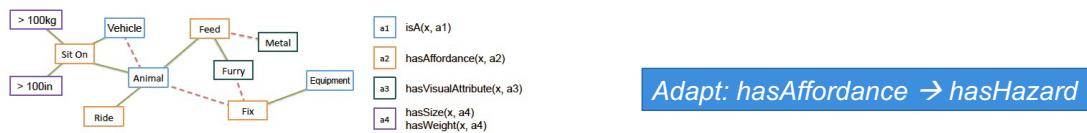


Fig. 7. Graphical illustration of the constructed KB. The nodes denote the entities (atomic formulae in MLN) illustrated on the right. The edges denote the attribute-attribute and attribute-affordance relations. The green solid edges indicate positive weights and the red dashed edges indicate negative weights.

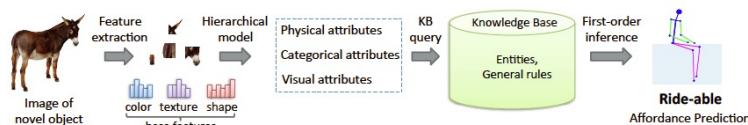


Fig. 9. The inference procedure of zero-shot affordance prediction. Given an image of a novel object, our model estimates the object attributes via a hierarchical model. These attributes serve as evidence for KB queries. We then employ first-order probabilistic inference to predict the affordances and to estimate human poses and human-object relative locations.



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 39

39

Deffordance

Detrimental affordance



a ladder is “climbable”



ladders are also “fall off -able”

For machines to reason about harms, they must be able to recognize hazards that endanger humans. We are extending for this purpose a combination of computer vision and knowledge engineering techniques that deduce object affordances (a ladder is “climbable”) to additionally deduce the detriments presented by objects (ladders are also “fall off -able”).



DISTRIBUTION A. Approved for public release: distribution unlimited.

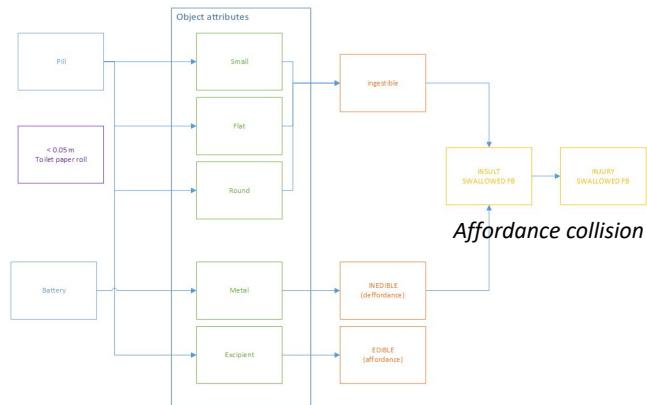
[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 40

40

Recognizing detriments

“ingestion of a foreign body”

- In a general sense, harms might be formulated as a product of the collision of an affordance and a dis-affordance (or “deffordance”).
 - Consider the diagnosis “*ingestion of a foreign body*.
 - That foreign body (e.g., button battery) has some physical size and shape properties (small, flat, round) that make it ingestible (affordance),
 - but in contrast to edible objects with those properties (e.g., pill), those objects with the material property of *inorganic* are inedible (deffordance).
 - This conclusion of detriment at the coincidence of *ingestible* but not *edible* is possible to reason to from object attributes.



80 YOM. SWALLOWED HEARING AID BATTERIES AFTER MISTOOK THE BATTERIES THE PILLS HE MEANT TO TAKE. DX: SWALLOWED FB

APL

DISTRIBUTION A. Approved for public release: distribution unlimited.

FOUNDATIONS | PERCEPTION | REASONING | 21 March 2022 | 41

Relationships of harm

Dangerousness is an inference



- Object O, with attributes $\{A_o\}$ that present hazards $\{H_o^A\}$
 - isDangerous To
 - Person P with attributes $\{A_p\}$ [propensity to $__$ / susceptibility] [intention/behavior] that impinge upon vulnerabilities $\{V_p^A\}$ **
 - (In Context C)
 - Because: A_o presents Hazard \rightarrow Vulnerability(H_o^A , V_p^A) via Mechanism of injury

C3 Object Balloon,
NR> material prop.
(rubber, so elastic)
(small, so swallowable)
presents C2/m1 5.1
mechanical Threat

KidA, a young child,
susceptible due to
(C: engaging in active
place L)

GAPS

 - Property of object
• presents hazards
 - Properties of individual
• that increase susceptibility

The beginnings of a notational depiction, linking to the classification scheme offered by WHO's International Classification of External Causes of Injury (ICECI)

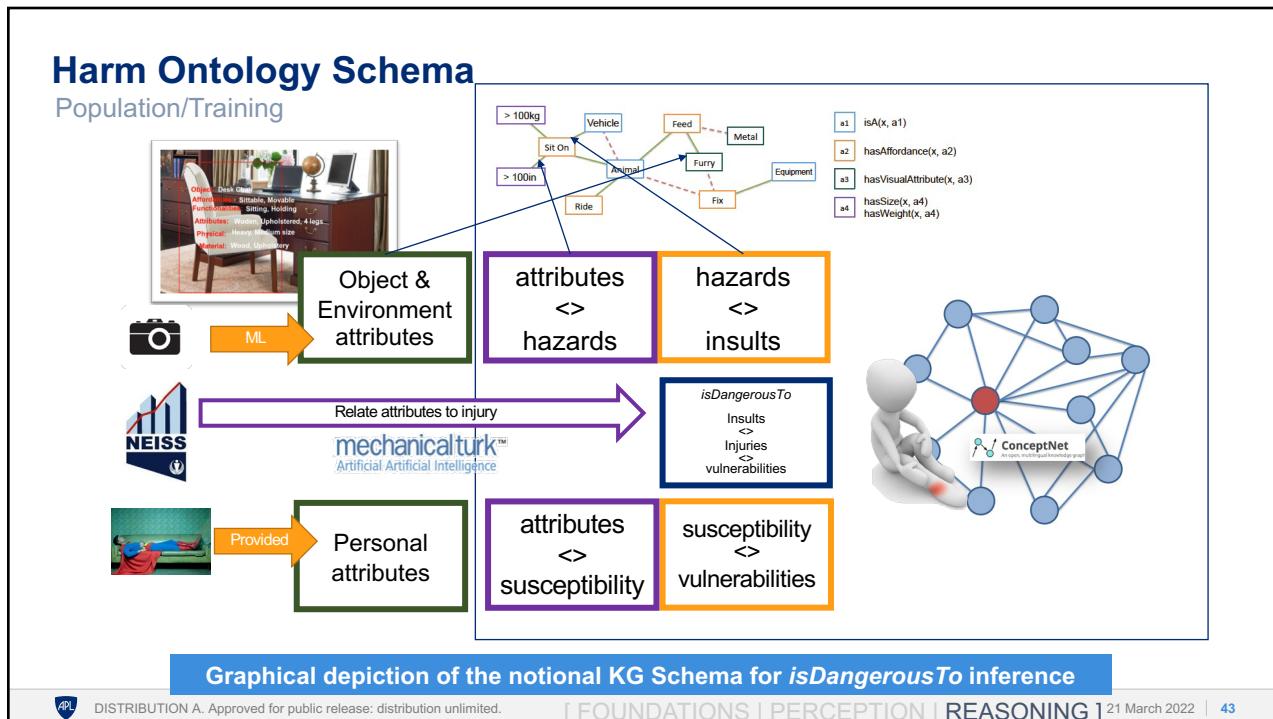
Formal logic

Dangerousness is not an intrinsic property of an entity, but rather an inference (*isDangerousTo*) relating an entity attribute to a vulnerability, and a human's particular susceptibility to be harmed by exposure to that attribute.

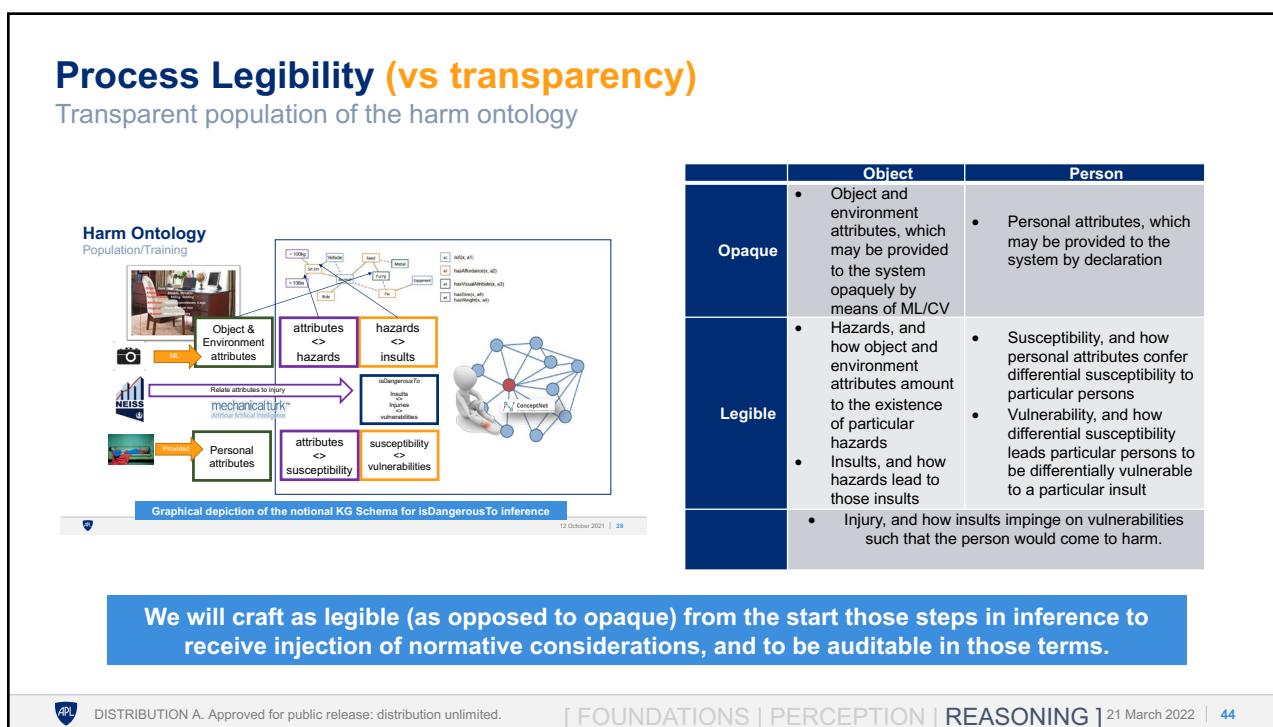
 APL DISTRIBUTION A. Approved for public release; distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 42

42



43



44

mTurk

Human annotation for populating ontology

Amazon Mechanical Turk (mTurk) is a crowdsourcing marketplace that makes it easier for individuals and businesses to outsource their processes and jobs to a distributed workforce who can perform these tasks virtually.



- For machines to generalize this nascent harms-reasoning capability to novel circumstances, we ask:
 - What are the attributes of each element within the epidemiological model of injury that allow harm to come to the host?
- Humans are naturally adept at responding to this kind of question.
 - Thus, we intend to construct microtasks wherein human participants perform structured decomposition of incident reports in NEISS to annotate with attributes the elements of product (object), person, and environment implicated in the injury.
- Importantly, with this reasoning so captured, novel inferences will be legible in terms of attributes, as compared to a black box approaches which are necessarily not so.



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING]

21 March 2022 | 45

45

Harm Ontology

Mechanical Turk task to populate relevant attributes

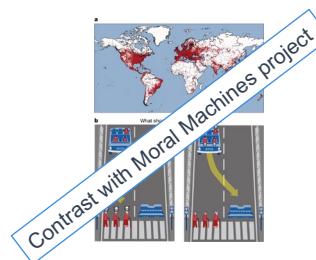
80 YOM. SWALLOWED HEARING AID BATTERIES AFTER MISTOOK THE BATTERIES THE PILLS HE MEANT TO TAKE. DX: SWALLOWED FB



- List the likely attributes of the PROD that enabled this to happen
 - Batteries resemble pills
 - PROD attributes:
 - Small, round, flat
 - affordance: swallowable
- if any, which attributes of the PERS make them particularly susceptible to the injury
 - Age degrades vision
 - Personal attributes:
 - AGE [NOT SEX, NOT RACE]
 - ETOH NOT CONSUMED
- if any, which attributes of the ENVI make them particularly susceptible to the injury
 - Age increases ambient abundance of pills and hearing aid batteries

Detrimental affordance (deffordance):
Ingested foreign body:

SWALLOWABLE (affordance)
BUT NOT
EDIBLE (affordance)



mechanicalturk
Artificial Artificial Intelligence



DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING]

21 March 2022 | 46

46

Harm Ontology

Leveraging emerging industry standards

DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 47

47

Summary

Knowledge Graph of Physical harms

The knowledge graph enables machine estimation of jeopardy, and supports the following functions:

- Recognizing detriments from attributes
- Balancing incommensurable harm types
- Parameterization of harms
- Generating value-driven behavior

Edge list
Results from ConceptNet 5.8
Source: Open Mind Common Sense contributors elper and borsos

Object	Relationship	Object	Weight
fire	— CapableOf —>	harm	0.0
too much of anything	— CapableOf —>	harm you	0.0

Documentation FAQ Chat Blog

Fig. 3. Object images and affordance labels. We illustrate 10 objects in our KB and their affordance labels. The x-axis lists the 14 affordances and the y-axis provides the 10 objects. The matrix shows that each object has many affordances, but each affordance is associated with FEW objects. This illustrates the concept of 'FEW affordances, compared to MANY defordances'.

DISTRIBUTION A. Approved for public release: distribution unlimited.

[FOUNDATIONS | PERCEPTION | REASONING] 21 March 2022 | 48

48

Enabling machines to reason about potential harms to humans

Outline

- Foundations
 - Ethics with respect to Artificial Agency
- Perception
 - Moral-scene assessment (recap)
- Reasoning
 - Harm ontology
- Upcoming Publications

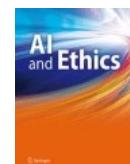
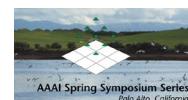


16 March 2022 | 49

49

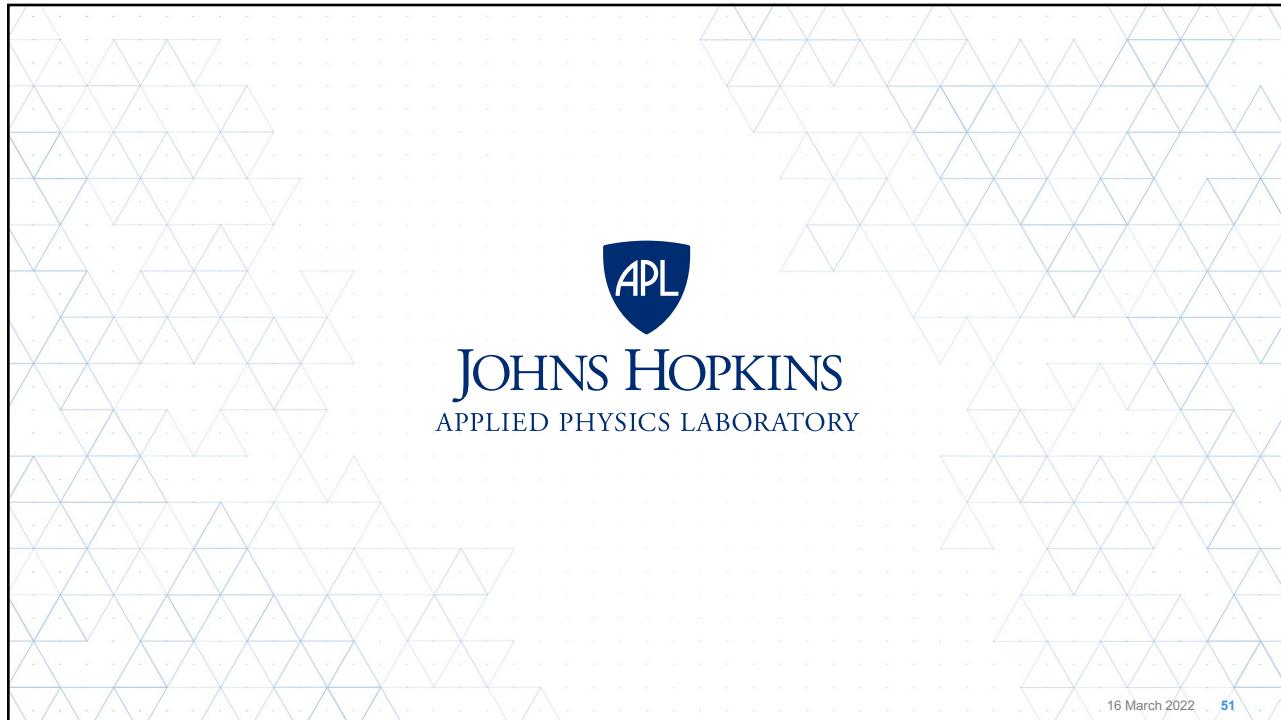
Upcoming publications

- **AAAI Spring Symposium (22 March 2022)**
 - Invited Talk in session *Putting AI in the Critical Loop: Assured Trust and Autonomy in Human-Machine Teams*
 - Follow-on book chapter: *Enabling machines to reason about potential harms to humans*
- **Frontiers in Physics: (14 May 2022)**
 - Interdisciplinary Approaches to the Structure and Performance of Interdependent Autonomous Human Machine Teams and Systems (A-HMT-S)
- **Entropy: Special issue (31 July 2022)**
 - Foundational matters in human interaction with machine agency
- **Springer: AI and Ethics**
 - Robots that *Do No Harm*



16 March 2022 | 50

50



51