

Endsem

Table of contents

- [Endsem](#)
- [Table of contents](#)
- [Robot Ethics](#)
- [Introduction to Robot Ethics by Patrick Lin](#)
 - [Overview and Significance](#)
 - [Historical Context and Cultural Impact](#)
 - [Robots in Society](#)
 - [Ethical and Social Issues](#)
 - [1. Safety and Errors](#)
 - [Critical questions raised:](#)
 - [2. Law and Ethics](#)
 - [Critical questions raised:](#)
 - [3. Social Impact](#)
 - [Critical questions raised:](#)
 - [Urgency and Proactivity in Ethics](#)
 - [Conclusion and Call to Action](#)
 - [References and Citations](#)
 - [Final Insight:](#)
- [Current Trends in Robotics: Technology and Ethics](#)
 - [Introduction and Significance](#)
 - [Definition of a Robot \(Section 2.1\)](#)
 - [Global Developments in Robotics \(Section 2.2\)](#)
 - [Industrial/Manufacturing Robots and Ethical Issues \(Section 2.3\)](#)
 - [Human-Robot Interaction in Healthcare, Surgery, and Rehabilitation \(Section 2.4\)](#)
 - [Robots as Co-inhabitants and Humanoid Robots \(Section 2.5\)](#)
 - [Socially Interactive Robots \(Section 2.6\)](#)
 - [Military Robots and Ethical Concerns \(Section 2.7\)](#)
 - [Conclusion \(Section 2.8\)](#)
 - [Notable Quotes:](#)
 - [Key References Mentioned:](#)
 - [Final Reflection:](#)
- [Roboethics: The Applied Ethics for a New Science](#)
 - [Introduction and Conceptual Clarification](#)
 - [Robotics as an Emerging Discipline \(Section 22.1\)](#)
 - [The Robotics Ideology \(Section 22.2\)](#)
 - [Robots and Moral Agency \(Section 22.3\)](#)
 - [Roboethics as a Work in Progress \(Section 22.4\)](#)
 - [Principles over Regulations: Military Robotics \(Section 22.5\)](#)
 - [Conclusion \(Section 22.6\)](#)
 - [Notable Quotes:](#)

- Key References Cited:
- Final Reflection:
- Robots In War
- Autonomous Driving Ethics: from Trolley Problem to Ethics of Risk
 - **1. Critique of the Trolley Problem and Transition to Risk Ethics**
 - **2. Limitations of Traditional Ethical Theories**
 - **3. Ethics of Risk: A Hybrid Framework**
 - **4. Application to the Trolley Problem**
 - **5. Social and Technical Implications**
 - **6. Unresolved Challenges and Future Directions**
 - **Conclusion**
- Just War and Robots' Killings
 - **1. Framing the Debate: Sparrow's Responsibility Trilemma**
 - **2. Rejecting the Trilemma: Engineering and Tolerance Levels**
 - **3. Ethics of Risk and the Precaution Thesis**
 - **4. Addressing Objections: Respect and Normalization**
 - **5. Practical Implications: Regulation Over Banning**
 - **6. Unresolved Tensions and Limitations**
 - **7. Synthesis with Just War Theory**
 - **8. Critical Examples and Analogies**
 - **9. Conclusion**
- Killer Robots by Robert Sparrow
 - **1. Core Argument: The Responsibility Trilemma**
 - **2. Ethical Foundations: *Jus in Bello* and Respect for Persons**
 - **3. The Child Soldier Analogy**
 - **4. Technological and Military Pressures**
 - **5. Unresolved Tensions and Counterarguments**
 - **6. Conclusion: Ethical and Policy Implications**
- Embedding values in AI
- How to Design AI for Social Good: Seven Essential Factors by Luciano Floridi, Josh Cowls, Thomas C. King, and Mariarosaria Taddeo, published in *Science and Engineering Ethics* (2020).
 - **Introduction: Framing AI for Social Good (AI4SG)**
 - **Methodology: How the Factors Were Identified**
 - **The Seven Essential Factors: Detailed Analysis**
 - **1. Falsifiability and Incremental Deployment**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**
 - **2. Safeguards Against the Manipulation of Predictors**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**

- **3. Receiver-Contextualised Intervention**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**
- **4. Receiver-Contextualised Explanation and Transparent Purposes**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**
- **5. Privacy Protection and Data Subject Consent**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**
- **6. Situational Fairness**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**
- **7. Human-Friendly Semanticisation**
 - **Explanation**
 - **Examples and References**
 - **Quotes**
 - **Best Practice**
 - **Analysis**
- **Balancing the Factors**
- **Conclusion and Future Directions**
- **Appendix: Representative Cases**
- **Embedding Values in Artificial Intelligence (AI) Systems**, published in *Minds and Machines* (2020).
 - **Introduction: Framing the Ethical Challenge in AI**
 - **Conceptualizing Values: A Normative Foundation**
 - **Embodied Values: A Triadic Distinction**
 - **AI Systems as Sociotechnical Systems: Five Building Blocks**
 - **Value Embedding in Technical Artifacts**
 - **Value Embedding in Institutions**
 - **Human Agents: Mediators of Value**
 - **Artificial Agents: Designed Autonomy**
 - **Technical Norms: Regulating AAs**
 - **System-Level Value Embedding**
 - **Conclusion: Practical Lessons**
- **Building Ethics into Artificial Intelligence** by Han Yu et al.,

- Introduction: Setting the Stage for Ethical AI
- Exploring Ethical Dilemmas: Understanding Human Preferences
 - GenEth: Expert-Driven Ethical Analysis
 - Moral Machine: Crowdsourcing Ethical Preferences
- Individual Ethical Decision Frameworks: Empowering Single Agents
 - MoralDM: Combining Rules and Analogies
 - BDI-Based Ethical Judgment
 - Game Theory and Machine Learning
 - CP-Nets for Preference Balancing
 - High-Level Action Language
 - Ethics Shaping in Reinforcement Learning
- Collective Ethical Decision Frameworks: Group Dynamics
 - Social Norms and Trust Networks
 - Human-Agent Collectives
 - Voting-Based System
- Ethics in Human-AI Interactions: Influencing Behavior
 - Persuasion Agents
 - Emotional Responses
- Discussions and Future Directions

Robot Ethics

Introduction to Robot Ethics by Patrick Lin

Overview and Significance

The chapter opens by contextualizing the current phase of robotics development as akin to the computer revolution, citing Bill Gates' notable assertion: "The emergence of the robotics industry is developing in much the same way that the computer business did 30 years ago" (Gates, 2007). Gates' comparison signifies that just as computers drastically reshaped society with both opportunities and challenges, robotics are expected to become ubiquitous, bringing complex social and ethical implications.

Historical Context and Cultural Impact

Lin draws attention to the historical fascination and apprehension society has had toward robots, emphasizing cultural references as far back as Homer's *Iliad* describing "golden servants" crafted by Hephaestus, to Leonardo da Vinci's mechanical knight, up to modern cultural representations such as *Metropolis*, *Blade Runner*, *Terminator*, and *I, Robot*. He suggests these cultural narratives reflect ongoing societal fears about robots' unpredictability and the potential dangers posed by emergent behavior or programming flaws.

Robots in Society

Robots are generally tasked with jobs humans find undesirable—dull, dirty, and dangerous (the "three Ds"). Lin lists extensive examples demonstrating the versatility of robots:

- **Labor and Services:** Factory robots in manufacturing, domestic robots like Roomba vacuums, and other service-oriented robots assisting with household tasks.
- **Military and Security:** Military robots named Predator, Reaper, and Big Dog, used for surveillance, bomb disposal, or combat. Civilian counterparts include police robots for security and home surveillance systems capable of dispensing pepper spray.
- **Research and Education:** Robots like NASA's Mars Exploration Rovers conducting space and oceanic research, or educational robots that interact directly with students.
- **Medical and Healthcare:** Robots such as the da Vinci Surgical System, therapeutic robots like PARO, and robotic nurses and pharmacists that help in clinical settings.
- **Personal Care and Companionship:** Robots assisting elderly and disabled individuals, exemplified by RI-MAN and CareBot, highlighting potential emotional bonds between humans and robots.
- **Environmental Management:** Robots addressing ecological challenges, such as cleaning oil spills, collecting toxic waste, or aiding after nuclear incidents.

Lin underscores that robotics is not static, but continuously evolving, potentially culminating in advanced nanobots or fully integrated biological-machine hybrids (cyborgs).

Ethical and Social Issues

Lin categorizes key ethical and social issues into three interrelated areas, each with explicit examples, historical context, and potential consequences:

1. Safety and Errors

Safety concerns are paramount because even minor flaws in programming can lead to fatal outcomes. Examples provided include:

- **Military Mishaps:** A U.S. military drone losing control and violating protected airspace in Washington D.C. (Bumiller, 2010), and a South African autonomous cannon malfunction killing friendly troops (Shachtman, 2007).
- **Civilian Fatalities:** The first robot-related fatality in a U.S. auto factory in 1979 (Kiska, 1983), illustrating real dangers beyond theoretical concerns.

Additionally, Lin highlights the risk of hacking, stressing that the characteristics making robots valuable—mobility, strength, and autonomy—could also be weaponized.

Critical questions raised:

- Can robots reliably differentiate between harmful and benign scenarios, like distinguishing weapons from benign objects?
- How should robots balance safety (e.g., kill-switches) against vulnerability to hacking?

2. Law and Ethics

Responsibility for harm caused by robots is legally ambiguous. Potential liable parties include developers, manufacturers, military commanders, or even the robots themselves, particularly as autonomy increases.

Lin also notes legal complications around privacy due to enhanced surveillance capabilities of robots, raising further concerns regarding invasive personal data collection by domestic robots connected to broader

networks.

Critical questions raised:

- How to encode ethical behavior in robots—should they follow deontological, consequentialist, or virtue ethics?
- Can or should robots be granted personhood, especially if integrated with human biology or consciousness?
- What ethical boundaries exist around substituting human relationships (companionship, caregiving) with robotic equivalents?

3. Social Impact

Robotics could significantly disrupt labor markets, paralleling the Industrial and Internet Revolutions. Lin acknowledges the common argument that automation frees humans for "higher-value" tasks, but emphasizes this is not universally comforting or viable, particularly for displaced workers who require immediate employment.

There's also concern about over-dependence on robotics eroding essential skills and creating societal fragility, illustrated historically by the Y2K computer crisis panic. Emotional relationships humans develop with robots present another dimension, with implications not yet fully understood, demonstrated by emotional attachment soldiers developed toward bomb-disposal robots (Singer, 2009a; Hsu, 2009).

Lin also identifies potential environmental harm, particularly through increased electronic waste and resource depletion, citing the ongoing e-waste crisis (O'Donoghue, 2010).

Critical questions raised:

- What societal structures and policies should we establish for handling job displacement?
- How should we handle increased dependency on robotic systems, and mitigate societal disruptions if robotic systems fail?
- Is emotional bonding with robots beneficial or potentially psychologically harmful?
- How will an expanded robotics industry impact global environmental sustainability?

Urgency and Proactivity in Ethics

Lin argues for proactive ethical consideration parallel to robotics development, highlighting historical delays in addressing ethical issues in technologies such as the Human Genome Project. He emphasizes the urgency of establishing ethical frameworks before widespread robot integration creates a "policy vacuum" (Moor, 1985).

Conclusion and Call to Action

Lin's overarching thesis emphasizes preparedness: as robotic capabilities rapidly expand, society must urgently confront these ethical dilemmas. Quoting Isaac Asimov, he stresses a proactive, science-fiction-inspired mindset:

"It is change, continuing change, inevitable change, that is the dominant factor in society today... our statesmen, our businessmen, our everyman must take on a science fictional way of thinking" (Asimov,

1978).

References and Citations

Throughout, Lin meticulously references authoritative sources:

- Historical/cultural references (Homer, da Vinci)
- Real-world incidents (Bumiller, Shachtman, Kiska)
- Ethical theories and frameworks (Arkin, Asimov)
- Current technological capabilities (Singer, Gates)
- Economic and environmental considerations (O'Donoghue, Rosenberg, Geipel)

These extensive citations ground Lin's arguments, providing scholarly depth and facilitating further exploration into specific subtopics.

Final Insight:

Patrick Lin systematically unpacks the profound, multifaceted implications of robotics advancement. His comprehensive overview serves not only as a guide to current ethical challenges but also as an imperative call for society to proactively engage these issues before they crystallize into intractable social problems. Lin's approach embodies a responsible, forward-looking ethical philosophy essential for guiding technological progress.

References directly from the provided text have been cited to maintain scholarly integrity.

Current Trends in Robotics: Technology and Ethics

Introduction and Significance

Bekey introduces the field of robotics as a "great technological success story" characterized by rapid advancements and increasing ubiquity in diverse fields—from healthcare and military to domestic chores and entertainment. Despite rapid technological growth, he notes, "the social and ethical implications of these new systems have been largely ignored," underscoring the urgency of ethical considerations parallel to technological development.

Definition of a Robot (Section 2.1)

Bekey first addresses the fundamental yet complex question: **"What is a robot?"** He clarifies that despite general fascination, a universally agreed-upon definition remains elusive. He proposes a working definition:

"A robot is a machine, situated in the world, that senses, thinks, and acts." (Bekey 2005)

This definition encapsulates several critical attributes:

- **Sensing:** Robots must gather data from their environment.
- **Thinking (processing):** Robots exhibit a degree of cognitive autonomy.

- **Acting (actuation):** Robots physically affect their environment through movement or force application.

Importantly, Bekey excludes purely virtual "software bots" or fully remote-controlled devices from his robot definition, emphasizing autonomy and environmental interaction as essential criteria.

Global Developments in Robotics (Section 2.2)

Bekey traces robotics evolution historically, noting an early American dominance that shifted toward Japan and Europe. He cites various companies like Unimation, Cincinnati Milacron (U.S.), Fujitsu, Panasonic, Kuka (Japan and Europe), highlighting shifts in global leadership over the decades.

He also identifies recent U.S. efforts to regain ground through government initiatives and roadmaps, contrasting this renewed activity with Europe's proactive ethical engagement, exemplified by the European Community's "Roboethics Roadmap" (*Veruggio 2006*).

Industrial/Manufacturing Robots and Ethical Issues (Section 2.3)

Bekey examines robotics' origins in manufacturing, referencing the introduction of "Unimate" (Engelberger, 1980) as pivotal. He recounts early fatal accidents, notably:

- A worker death at Ford's Michigan plant (1979).
- A robot killing a Japanese worker during maintenance (1981).

These incidents highlighted crucial ethical concerns about human-robot workplace interaction, prompting safety barriers. Ethical issues arising here include:

1. **Fear of Replacement:** Workers fearing job loss due to automation. Bekey advocates responsible management strategies—worker inclusion in planning, training for new roles—to mitigate anxiety.
2. **Dehumanization of Work:** Repetitive tasks given to robots may cause workers to feel inferior, sparking resentment reminiscent of 19th-century Luddites who opposed mechanization. Ethical management involves assigning tasks to humans leveraging their unique cognitive capabilities.
3. **Human-Robot Cooperative Work:** Recent "cobot" (collaborative robot) development aims to blend robot precision with human decision-making, significantly reducing risk and promoting safer human-robot interactions (*Gillespie et al., 2001*). Despite benefits, such cooperation may unintentionally reduce vital human-to-human workplace interactions, raising ethical considerations needing proactive resolution.

Human-Robot Interaction in Healthcare, Surgery, and Rehabilitation (Section 2.4)

Healthcare robotics, including nursing and surgery, exemplify rapidly expanding human-robot interactions. Bekey details developments like the robotic assistant "HelpMate," capable of independently navigating hospital environments, and "Pearl," assisting elderly patients by reminding medications and offering companionship (*Montemerlo et al., 2002*).

Key ethical concerns include:

- **Emotional attachments:** Patients becoming overly dependent on robots.

- **Robotic limitations:** Robots may inadequately handle emotional patient responses or complex ethical decisions (e.g., medication refusals).

In robotic surgery, Bekey highlights the "da Vinci Surgical System," emphasizing its current status as teleoperated (human-controlled remotely), yet raising crucial ethical questions anticipating future autonomous robotic surgeons:

- "If complications arise, who bears responsibility—the designer, manufacturer, surgeon, hospital, or insurer?"
- Determining ethically acceptable levels of risk in robotic surgeries and appropriate accountability mechanisms.

Robots as Co-inhabitants and Humanoid Robots (Section 2.5)

Robots like "Roomba" vacuum cleaners demonstrate domestic robotic success, paving the way for more advanced "humanoid robots" that share living spaces with humans. Humanoids such as "Wakamaru" (Japan) and "Nao" (France) exemplify sophisticated robotic cohabitants capable of complex human interactions, including gesture recognition and responsive communication.

Bekey addresses ethical questions like:

- Privacy invasions by robots within homes.
- Potential misuse, e.g., robots programmed to engage in unethical behaviors.
- Whether robots deserve rights or respectful treatment analogous to humans.
- Managing emotional interactions or ethical responses to robot malfunctions or perceived misconduct.

Socially Interactive Robots (Section 2.6)

This broader category encompasses robots designed explicitly for social engagement. Bekey cites extensive ongoing research into robot swarms, group behaviors, and robot-to-robot interactions, potentially leading to complex societies with robots exhibiting unique personalities and advanced communication.

He also discusses research into robot emotional expressivity, referencing MIT's "Kismet" robot (Breazeal, 2002), designed to exhibit human-like emotions, significantly influencing human interaction and raising ethical questions about anthropomorphism and appropriate human reactions to robot-expressed emotions.

Military Robots and Ethical Concerns (Section 2.7)

Military robots, used extensively for explosive disposal and combat operations, generate critical ethical concerns. Bekey discusses hypothetical scenarios illustrating ethical dilemmas, such as:

1. A robot discovering noncombatant children in a targeted building, conflicting with programmed engagement rules.
2. Autonomous drone self-defense potentially harming humans due to time-critical decision-making constraints.

Citing research by Arkin (2009), Bekey emphasizes current inadequacies in existing ethical frameworks (laws of war, rules of engagement), which robots might struggle to interpret correctly. He poses vital ethical questions, including:

- Human rights violations by autonomous robots.
- Responsibility for collateral property damage or casualties.
- Risks of lowered barriers to war participation due to perceived reduced casualties from robot deployment.
- Risks of robotic technologies proliferating internationally, complicating global security.

Conclusion (Section 2.8)

Bekey concludes by reiterating the gap between rapid technological advancements and lagging ethical deliberation within robotics communities. He emphasizes the need for ongoing ethical reflection, proactive governance, and comprehensive understanding of the social implications accompanying robotics integration.

Notable Quotes:

- "Only during the past decade have we seen the emergence of the field of 'robot ethics'... with most efforts in Europe, Asia, and the United States."
- "Management has an ethical responsibility to allow humans to work in tasks that do not demean them, but rather take advantage of their superior cognitive abilities."
- "The risks of using a robot surgeon... must be lower than those encountered with human surgeons."
- "Humans have a tendency to anthropomorphize robots, and any display of emotions (real or artificial) by the robot could lead to unacceptable (or unethical) behaviors by humans."

Key References Mentioned:

- Bekey, G.A. (2005). *Autonomous Robots*.
- Breazeal, C. (2002). *Designing Sociable Robots*.
- Arkin, R.C. (2009). *Governing Lethal Behavior in Autonomous Robots*.
- Singer, P. (2009). *Wired for War*.
- Veruggio, G. (2006). EURON Roboethics Roadmap.
- Gillespie, R.B., Colgate, J.E., Peshkin, M.A. (2001). Framework for cobot control.

Final Reflection:

Bekey provides a thorough exploration of current robotics developments and their accompanying ethical challenges, urging the field to balance innovation with proactive ethical engagement. His comprehensive survey illustrates the depth and breadth of robotics' potential impacts on society, emphasizing that ethical deliberation must occur concurrently with technological advancement to ensure beneficial integration into human environments.

Here's a detailed, in-depth analysis of **Chapter 22: "Roboethics: The Applied Ethics for a New Science"** by **Gianmarco Veruggio and Keith Abney** from the book *Robot Ethics: The Ethical and Social Implications of Robotics*, thoroughly explaining key concepts, including quotes, references, examples, and implications without summarizing superficially:

Roboethics: The Applied Ethics for a New Science

Introduction and Conceptual Clarification

The chapter begins by recognizing robotics as a relatively new scientific discipline that raises complex ethical issues. "Roboethics" is introduced as a field specifically addressing ethical considerations around robotics and its social implications. The authors emphasize that "robot ethics" can have **three distinct meanings**:

1. **Applied Ethics**: Exploring ethical and social implications of robotic technology in human society.
2. **Programmed Ethics**: Ethical codes embedded in robots by human programmers.
3. **Autonomous Moral Agency**: The hypothetical scenario where robots possess self-conscious ethical reasoning capabilities, becoming full moral agents.

Veruggio introduces the term "**roboethics**" specifically for the first meaning, defining it as:

"An applied ethics whose objective is to develop scientific, cultural, and technical tools that can be shared by different social groups and beliefs" (*Veruggio, 2007*).

This human-centered perspective places ethical responsibility on humans—researchers, designers, and users—not robots themselves.

Robotics as an Emerging Discipline (Section 22.1)

The authors delve into robotics as an evolving branch of engineering, defined succinctly as:

"A robot is a machine, situated in the world, that senses, thinks, and acts" (*Bekey, 2005*).

Robotics, as explained, signifies the "Third Industrial Revolution," with machines increasingly capable of autonomous decisions and interactions. Such autonomy brings profound ethical considerations. Robotics is described as inherently interdisciplinary, necessitating philosophical, psychological, sociological, and legal insights.

The Robotics Ideology (Section 22.2)

A key theme here is the role of ideology—culturally ingrained myths and misconceptions—in shaping public perception of robots. Popular fears, like the trope of a robotic uprising ("Rebellion of the Automata"), are cited as examples of ideology rather than scientific realism. Veruggio and Abney criticize these myths as:

- Unrealistic, driven by irrational fears or guilt related to historical slavery.
- Misleading public expectations away from actual robotic capabilities and immediate ethical concerns.

They reference iconic fictional robots like **HAL 9000** from Arthur C. Clarke's *2001: A Space Odyssey* and **replicants** from Philip K. Dick's *Do Androids Dream of Electric Sheep?*, illustrating how unrealistic fictional scenarios have distorted public understanding and expectations of robots.

Contrasting Western and Japanese cultures, the authors observe that Japanese Shinto traditions, which blur animate-inanimate boundaries, lead to more positive attitudes toward robots compared to Western anxieties.

The "Pinocchio Syndrome"—the erroneous belief that robots could literally evolve into humans—is criticized as a fallacy conflating functional equivalence with actual biological or ontological identity. Robots might achieve symbolic reasoning abilities but could never literally become human beings.

Robots and Moral Agency (Section 22.3)

Central to roboethics is the question of whether robots could ever achieve moral agency—the capacity for ethical reasoning, self-consciousness, and responsibility. The authors thoroughly explore potential criteria for robotic moral agency, discussing Kant's "transcendental unity of apperception" (TUA), free will, symbolic reasoning, and embodiment theories (Embodied Cognition).

- Kant's concept of TUA suggests robots must achieve unified, self-aware consciousness—something currently beyond robotic capabilities.
- Free will, as posited by philosopher José Galván, is identified as critical to genuine moral agency:

"Free will is a condition of man, which transcends time and space... [It] cannot be imitated by a machine" (*Galván, 2004*).

- Embodied Cognition (EC), advocated by philosophers like Rodney Brooks and Lakoff & Johnson, emphasizes physical embodiment as critical to consciousness and moral agency, challenging the idea that purely computational minds could achieve genuine consciousness.

The authors acknowledge ongoing philosophical debate (e.g., Searle's "Chinese Room" argument and Churchlands' neurophilosophical positions) without prematurely settling these profound issues.

Roboethics as a Work in Progress (Section 22.4)

The practical implications of roboethics extend to urgent contemporary matters, highlighting that ethical guidelines must develop in parallel with technological innovation. The Roboethics Roadmap initiated by Veruggio after the First International Symposium on Roboethics (2004) exemplifies a proactive, interdisciplinary approach—integrating scientists, ethicists, and policymakers—to ensure robotics progresses ethically and safely.

Important practical considerations include:

- Implementing safety standards for autonomous robots.
- Defining clear legal frameworks governing robot mobility, accountability, and liability.
- Ensuring ethical deliberation, rather than profit-driven market forces alone, guides robotic development.

The authors underscore that ethical decision-making in robotics cannot be neutral; abstaining from regulation inherently favors powerful economic interests over societal welfare:

"To avoid regulation is itself a choice... Abstention ultimately ends up favoring the strongest" (*Coiffet, 2004*).

Principles over Regulations: Military Robotics (Section 22.5)

Military robotics exemplifies critical ethical challenges needing immediate attention. The authors argue ethical principles must precede technical regulations. Military robots, promoted as advantageous due to performing tasks described as **dull, dirty, dangerous, and dispassionate**, nonetheless raise severe ethical concerns:

- Reliability and precision in distinguishing combatants from civilians.

- Autonomy potentially shifting accountability away from human operators or commanders.
- Possibility of lowering the threshold for warfare, given reduced human casualties on the deploying side.

Historical analogies, such as the Saint Petersburg Declaration (1868) against certain munitions, illustrate past attempts—and failures—to limit warfare's cruelty through ethical agreements. The authors caution that optimistic claims about ethical robotic soldiers adhering perfectly to international humanitarian laws remain unrealistic given current technological limitations:

"Until fully autonomous robots demonstrate (in realistic simulations) that they are no more likely to commit war crimes than human soldiers, it seems immoral to deploy them."

Conclusion (Section 22.6)

In concluding, the authors reinforce that roboethics requires collaborative, multidisciplinary dialogue among scientists, ethicists, policymakers, and the public. They advocate for:

- Dispelling popular misconceptions through informed public debate.
- Developing cross-cultural ethical frameworks adaptable to international laws.
- Prioritizing human-centered ethical considerations over market-driven technological advances.

Ultimately, the authors warn that neglecting roboethics risks severe societal harm:

"It is crucial to tackle not the mythical worries due to ideologies... but the real issues facing robotics in the larger society—before it's too late."

Notable Quotes:

- "The explicit aim... is to develop autonomous robots that substitute for human soldiers... untiring and near-invincible robotic soldiers."
 - "Popular misconceptions... largely stem not from its being a new scientific discipline, but from its status as an ideology."
 - "Perhaps the worries over the so-called rebelling automata are because we think of them... as human slaves."
 - "Before discussing 'how,' we should decide 'if' a fully autonomous robot can be allowed to kill a human."
-

Key References Cited:

- Bekey, G. (2005). *Autonomous Robots*.
 - Coiffet, P. (2004). Speech on humanist development of robotics.
 - Galván, J. (2004). Technoethics and free will.
 - Kant, I. ([1781/1787] 1997). *Critique of Pure Reason*.
 - Searle, J. (1984). "Minds, Brains, and Science".
 - Veruggio, G. (2007). EURON Roboethics Roadmap.
 - Warwick, K. (2002). *I, Cyborg*.
-

Final Reflection:

This chapter compellingly demonstrates that ethical inquiry must proceed alongside technological advancements in robotics, avoiding sensational myths and prioritizing human welfare and global dialogue. Roboethics is not merely philosophical but an urgent, practical imperative as robotics increasingly integrates into daily life and societal infrastructures.

Robots In War

Autonomous Driving Ethics: from Trolley Problem to Ethics of Risk

1. Critique of the Trolley Problem and Transition to Risk Ethics

The paper begins by dismantling the trolley problem's relevance to autonomous vehicles (AVs), arguing that its **deterministic, binary framework** fails to capture real-world complexities. The authors highlight three critical shortcomings:

- **Certainty of Outcomes:** Trolley scenarios assume fatalities are inevitable, whereas real driving involves probabilistic risks. For example, collisions depend on sensor accuracy, prediction errors, and occluded objects (Fig. 4), making absolute certainty unrealistic.
- **Binary Choices vs. Continuous Solutions:** Unlike the trolley's two-track dilemma, AVs navigate a "continuous solution space for trajectories" (p. 3), requiring nuanced risk assessments across countless potential paths.
- **Lack of Context:** Trolley problems omit critical prior information (e.g., who caused the risk), which is ethically essential. Nyholm & Smids (2016) and Kauppinen (2020) emphasize that moral responsibility hinges on contextual factors like fault or intent, which the trolley framework ignores.

The authors pivot to **ethics of risk**, which evaluates actions under uncertainty. This shift aligns with Bonnefon et al. (2019), who reframe AV decisions as probabilistic risk distributions rather than life-and-death trade-offs. For instance, an AV adjusting its lateral position (Fig. 2) redistributes risk between cyclists and passengers based on collision probabilities and harm severity.

2. Limitations of Traditional Ethical Theories

The paper systematically critiques three ethical frameworks for AVs:

A. Deontology

Rule-based systems like Kant's Categorical Imperative or Asimov's Three Laws are deemed **too rigid** for dynamic environments. While Gerdes & Thornton (2015) propose hierarchical rules (e.g., prioritizing human safety over obedience), the authors argue such approaches:

- Fail to account for context (e.g., swerving to avoid a pedestrian might endanger others).
- Create conflicts between rules, leading to "dangerous behaviors" (p. 6) if rigidly enforced.
- Lack universality, as corner cases require endless rule additions.

B. Utilitarianism

Though utilitarian cost functions (minimizing total harm) are technically feasible and socially preferred (Bonnenfon et al., 2016), they face **moral and legal challenges**:

- Sacrificing individuals for collective benefit violates human dignity, as per the German Ethics Commission (2017).
- Transparency issues arise when AVs prioritize "invisible" statistical lives over identifiable passengers (Hubner & White, 2018).
- Ignores fairness, as noted in Keeling’s (2018) critique of Leben’s (2017) Rawlsian algorithm, which disproportionately favors worst-off individuals.

C. Virtue Ethics

Machine learning (ML)-based approaches, like imitation learning (Bansal et al., 2018), aim to encode "virtues" such as prudence. However, the authors identify **critical flaws**:

- ML models lack explainability, making accountability impossible (Berberich & Diebold, 2018).
- Training data reflect common behaviors, not ethical ideals (Etzioni & Etzioni, 2017). For example, AVs might mimic aggressive human driving rather than virtuous caution.

3. Ethics of Risk: A Hybrid Framework

The authors propose a **risk-based model** integrating three decision principles:

1. **Bayesian Principle:** Minimize total risk ($J_B = \sum R_i$).
2. **Equality Principle:** Reduce risk disparities ($J_E = \sum |R_i - R_j|$).
3. **Maximin Principle:** Mitigate worst-case harm ($J_M = \max(H_i)$).

These are combined into a weighted cost function:

[$J_{\text{total}} = (w_B J_B + w_E J_E + w_M J_M) \gamma^t$]

where (γ^t) discounts future risks.

Key Innovations:

- **Probabilistic Harm Quantification:** Risk ($R = p(u)H(u)$) incorporates collision probability (p) and harm severity (H), estimated via kinetic energy models (Sobhani et al., 2011). This avoids subjective valuations (e.g., monetizing lives) and focuses on physical metrics like speed and mass.
- **Fairness Through Hybridization:** By balancing total risk reduction (Bayesian), equity (Equality), and worst-case avoidance (Maximin), the framework addresses Keeling’s (2018) critiques. For example:
 - In Fig. 5, the Equality Principle alone would favor two certain deaths over one minor risk, but combining it with Bayesian weights ($w_B > 0$) corrects this.
 - In Fig. 6, the Maximin Principle ignores cumulative harm to many, but adding Bayesian terms prioritizes collective safety.

Limitations:

- **Subjective Weightings:** The choice of (w_B, w_E, w_M) is unresolved. Manufacturers might prioritize passenger safety (high (w_M)), contradicting the German Ethics Commission’s rejection of "sacrifice" logic.

- **Responsibility Metrics:** While the framework discounts risks temporally ((γ^t)), it doesn't penalize negligent road users (e.g., jaywalkers). Kauppinen (2020) argues culpability must influence risk distribution, but the paper defers this to future work.

4. Application to the Trolley Problem

The authors demonstrate their framework's flexibility by applying it to the trolley dilemma:

- **Standard Scenario:** Killing one vs. five. With $(w_B > 0)$, Bayesian dominance $((J_B = 1)$ vs. (5)) yields the utilitarian outcome.
- **Modified Scenario:** Five suffer minor harm $((H = 0.2))$ vs. one killed. Here, Maximin $((J_M = 0.2))$ and Equality $((J_E = 1))$ override Bayesian $((J_B = 1))$, favoring the five.

This illustrates the framework's adaptability but also its dependence on weightings. The authors concede that without predefined weights, outcomes can be ambiguous, stating:

"the question of the weighting factors [...] cannot be answered separately" (p. 17).

5. Social and Technical Implications

- **Mandatory vs. Personal Ethics:** The framework accommodates both approaches. Weights could be standardized (mandatory) or user-adjusted (personal "ethical knob" à la Contissa et al., 2017). However, the authors warn that personal settings might increase insurance costs or societal risk.
- **Transparency and Regulation:** Explicit mathematical models enable post-accident audits, addressing the "black box" critique of ML-driven systems. Regulators could mandate transparency in (w) values to ensure compliance with ethical guidelines.
- **Cultural Variability:** The paper acknowledges but does not resolve how cultural differences in risk tolerance (e.g., individualistic vs. collectivist societies) might influence weightings.

6. Unresolved Challenges and Future Directions

- **Quantifying Responsibility:** The framework doesn't penalize at-fault road users. Future work might assign risk discounts based on culpability (e.g., red-light violators).
- **Dynamic Weightings:** Context-dependent (w) adjustments (e.g., higher (w_M) in school zones) could enhance fairness.
- **Legal Integration:** The German Ethics Commission's guidelines (2017) reject utilitarian sacrifices but permit probabilistic risk minimization. Aligning the framework with such policies requires legal-philosophical reconciliation.

Conclusion

Geisslinger et al. (2021) provide a groundbreaking shift from abstract trolley dilemmas to actionable risk ethics for AVs. Their hybrid framework balances competing ethical principles and operationalizes them into tractable mathematics. However, the reliance on subjective weightings and unresolved responsibility metrics leaves critical gaps. Future research must address these to ensure AVs navigate not just roads, but moral landscapes, with rigor and fairness. As the authors conclude:

"The question of what constitutes a fair distribution of risk [...] should be at the center of future research" (p. 20).

Just War and Robots' Killings

1. Framing the Debate: Sparrow's Responsibility Trilemma

The paper opens by addressing Rob Sparrow's **responsibility trilemma**, a central ethical challenge to lethal autonomous weapons systems (LAWS). Sparrow argues that if a robot commits a war crime (e.g., targeting surrendering soldiers), there is no morally responsible agent:

- **Designers** are absolved because autonomous systems are designed to act unpredictably.
- **Commanders** cannot foresee autonomous actions, akin to artillery operators who are not blamed for misuse by soldiers.
- **Robots themselves** lack moral agency, as they cannot suffer or comprehend responsibility (pp. 304–305).

This creates a **responsibility gap**, rendering LAWS impermissible under Just War theory's requirement for accountability. Sparrow analogizes LAWS to child soldiers in "grey areas" of autonomy, where no party bears responsibility for atrocities (p. 305).

2. Rejecting the Trilemma: Engineering and Tolerance Levels

Simpson and Müller counter by introducing **tolerance levels**, a concept borrowed from engineering ethics. Tolerance levels define the reliability a system must achieve, integrating technical and moral considerations:

- **Example:** A bridge's collapse due to a 300-year flood event falls outside its tolerance level; no one is blameworthy. Conversely, failure due to poor design or misuse (e.g., overloading) assigns responsibility to engineers or users (pp. 307–309).
- **Application to LAWS:** If a robot operates within its tolerance level (e.g., targeting combatants with 99% accuracy), designers, commanders, or regulators are responsible for malfunctions. If it operates outside (e.g., unforeseeable sensor failure in extreme conditions), no blame accrues (pp. 310–311).

This framework rejects Sparrow's trilemma by distributing responsibility across a **chain of human actors**, similar to liability in civilian engineering projects.

3. Ethics of Risk and the Precaution Thesis

The authors pivot to the **ethics of risk imposition**, addressing whether LAWS can fairly redistribute risk. Key arguments include:

- **Risk Baselines:** The moral permissibility of LAWS hinges on whether they reduce non-combatant risk compared to human soldiers ($r_2 < r_1$). For example, if LAWS reduce collateral damage by 50%, their deployment is justified despite unequal risk distribution (pp. 314–315).
- **Policy G Analogy:** Reducing risk unequally (e.g., eliminating risk for soldiers while slightly reducing it for non-combatants) is acceptable if no egalitarian alternative exists. This aligns with Rawlsian prioritarianism, prioritizing absolute risk reduction over perfect equality (p. 315).

The **Precaution Thesis**—requiring precautions to minimize harm to each individual—is satisfied if LAWS meet two conditions:

1. ($r_2 < r_1$) (overall risk reduction).
2. (r_2) is minimized to the “technologically feasible” limit (p. 316).

4. Addressing Objections: Respect and Normalization

The paper confronts two objections:

- **Respect for Victims:** Sparrow and Nagel argue that LAWS violate interpersonal respect by depersonalizing killing. Simpson and Müller respond with **technological normalization**, comparing LAWS to artillery or drones:

“Familiarity with the weapon has rendered it no less disrespectful as a tool for killing than face-to-face combat” (p. 320).

Historical precedents (e.g., artillery) show that societal acceptance evolves, dissolving initial moral revulsion.

- **Psychological Aversion:** The “atavistic horror” of being killed by robots is likened to tax aversion—irrational but transient. Over time, LAWS would become routine, much like self-checkout systems (p. 318).

5. Practical Implications: Regulation Over Banning

The authors conclude with a pragmatic stance: **regulate, don’t ban**. Key regulatory imperatives include:

- Setting **strict tolerance levels** for LAWS to ensure compliance with Just War principles (e.g., discrimination, proportionality).
- Licensing and testing regimes to enforce accountability (p. 320).
- International governance to prevent misuse by “badly ordered societies” (p. 320).

6. Unresolved Tensions and Limitations

- **Threshold Objection:** LAWS might lower the threshold for war by making conflict less politically costly (e.g., reduced soldier casualties). This risks increasing unjust wars (p. 320).
- **Democratic Accountability:** Centralized control of LAWS could bypass public consent, undermining democratic oversight (p. 320).
- **Liability in Asymmetric Conflicts:** The paper assumes LAWS will be used by “well-ordered societies,” but real-world deployment (e.g., by non-state actors) complicates accountability.

7. Synthesis with Just War Theory

The authors align their argument with **revisionist Just War theory**, rejecting the “moral equality of combatants” and emphasizing liability based on threat contribution (pp. 310–311). They sidestep debates over combatant liability by focusing on risk redistribution, asserting that LAWS can satisfy *jus in bello* principles if calibrated to minimize non-combatant harm.

8. Critical Examples and Analogies

- **Mefloquine Case:** Pharmaceutical companies are not blamed for rare side effects if the drug meets regulatory tolerance levels (p. 309). Similarly, LAWS operators are blameless for statistically inevitable malfunctions.
- **Dam Construction:** A dam reducing flood risk, despite introducing a new collapse risk, is justified if ($r_2 < r_1$) (p. 314).

9. Conclusion

Simpson and Müller provide a robust defense of LAWS within Just War theory, reframing responsibility through engineering ethics and risk analysis. However, their reliance on **technological optimism** (e.g., assuming ($r_2 < r_1$)) and societal normalization leaves room for critique. As they concede:

“The deep and difficult question [is] whether it is likely that (r_2) will be greater or lesser than (r_1)” (p. 316).

The paper’s strength lies in its integration of moral philosophy with practical engineering standards, offering a blueprint for ethical LAWS deployment—if regulators enforce stringent tolerance levels and prioritize non-combatant safety.

Killer Robots by Robert Sparrow

1. Core Argument: The Responsibility Trilemma

Sparrow’s central thesis revolves around the **responsibility gap** inherent in deploying lethal autonomous weapon systems (AWS). He argues that AWS’s autonomy precludes meaningful accountability for war crimes, violating the ethical foundations of *jus in bello*. The trilemma posits three potential loci of responsibility—programmers, commanders, and machines—and systematically dismantles each:

1. Programmers:

- **Argument:** Programmers cannot foresee all actions of a learning, adaptive AWS. Autonomy breaks the causal chain between programming and outcomes.
- **Example:** A robot that evolves beyond its initial code to target surrendering soldiers (p. 66) cannot have its designers blamed, as its decisions reflect “internal states” (beliefs, desires) shaped post-deployment.
- **Limitation:** Assumes programmers cannot implement fail-safes or ethical constraints. Contemporary debates on “value alignment” in AI challenge this (e.g., embedding Asimov’s laws).

2. Commanding Officers:

- **Argument:** Commanders are likened to artillery operators, responsible for deployment but not specific targeting. However, AWS’s unpredictability makes this analogy flawed.
- **Quotation:** “The more autonomous the systems are, the larger this risk looms. At some point, it will no longer be fair to hold the Commanding Officer responsible” (p. 70).
- **Counterpoint:** Military doctrine already holds commanders accountable for collateral damage from “dumb” weapons (e.g., artillery). Why not extend this to AWS? Sparrow dismisses this by emphasizing AWS’s unique decision-making capacity (p. 70).

3. Machines Themselves:

- **Argument:** Machines lack moral agency. Punishing them is nonsensical, as they cannot "suffer" or comprehend guilt (p. 72).
- **Analogy:** Assigning blame to a machine is as absurd as prosecuting a malfunctioning toaster (p. 71).
- **Philosophical Challenge:** Sparrow presupposes a Kantian view of moral responsibility requiring consciousness. Proponents of "artificial moral agents" (Floridi & Sanders, 2001) argue responsibility could hinge on functional behavior, not sentience.

2. Ethical Foundations: *Jus in Bello* and Respect for Persons

Sparrow grounds his argument in Just War Theory's requirement for accountability:

- **Key Principle:** "The least we owe our enemies is allowing that their lives are of sufficient worth that someone should accept responsibility for their deaths" (p. 67).
- **Consequences:** Unaccountable AWS reduce warfare to "extermination," akin to using landmines or WMDs (p. 67).

Strengths:

- Highlights the dehumanizing effect of detached, automated killing.
- Aligns with Nagel's (1972) emphasis on interpersonal respect in war: Combatants deserve to know *who* decided their fate and *why*.

Weaknesses:

- **Technological Determinism:** Assumes AWS will inherently lack oversight. Yet, hybrid systems (e.g., human-AI collaboration) could retain accountability (e.g., requiring final human approval for strikes).
- **Historical Precedent:** Drones already operate with human oversight. Sparrow acknowledges this but argues autonomy trends will eliminate "the human in the loop" (p. 68–69).

3. The Child Soldier Analogy

Sparrow's most provocative comparison links AWS to child soldiers:

- **Shared Trait:** Both occupy a "grey area" of partial autonomy. Children lack full moral agency; AWS lack moral personhood.
- **Ethical Dilemma:** "No one is in control. If civilians are killed, they are killed senselessly without anyone being responsible" (p. 73).
- **Implication:** Using AWS mirrors the moral bankruptcy of deploying child armies (e.g., Liberia, Angola).

Critique:

- **False Equivalence:** Child soldiers are victims coerced into violence, whereas AWS are tools designed by humans. The analogy conflates *moral patients* (children) with *amoral instruments* (robots).
- **Policy Response:** International law bans child soldiers via the Optional Protocol to the CRC. A similar ban on AWS (e.g., proposed UN Treaty) could resolve Sparrow's trilemma.

4. Technological and Military Pressures

Sparrow identifies systemic drivers pushing toward AWS deployment:

1. **Tempo of Battle:** AI's speed outperforms human decision-making in air combat (p. 68).
2. **Cost-Efficiency:** AWS reduce soldier casualties and financial burdens (p. 64).
3. **Survivability:** Removing human operators mitigates communication vulnerabilities (p. 69).

Contradiction: While AWS promise precision (reducing collateral damage), their autonomy risks *increased* unpredictability. Sparrow cites the LOCAAS missile system, which autonomously selects warhead configurations (p. 64), as a harbinger of ethical chaos.

5. Unresolved Tensions and Counterarguments

- **Regulation vs. Ban:** Sparrow dismisses regulation ("Human oversight will eventually be eliminated," p. 68) but ignores frameworks like the EU's AI Act, which mandates human control over high-risk AI.
- **Moral Progress:** If AWS reduce civilian casualties compared to human soldiers, their use might be *more* ethical despite accountability gaps. Sparrow's focus on responsibility overlooks consequentialist gains.
- **Moral Agency Evolution:** Advances in AI consciousness (e.g., artificial general intelligence) could render AWS "moral persons." Sparrow rejects this (p. 72) but does not engage with futurists like Kurzweil (1999).

6. Conclusion: Ethical and Policy Implications

Sparrow's analysis remains a cornerstone in AWS ethics, compellingly arguing that autonomy erodes accountability. However, his conclusions face challenges:

- **Techno-Optimism:** Engineers like Arkin (2009) propose "ethical governors" to constrain AWS behavior, potentially resolving responsibility gaps.
- **Legal Precedent:** The Ottawa Treaty banned landmines; a similar ban on AWS could preempt Sparrow's dystopia.
- **Philosophical Evolution:** Debates on AI moral agency (e.g., machine consciousness) may redefine responsibility paradigms.

Ultimately, Sparrow's warning—that AWS risk rendering war "unfair either to potential casualties [...] or to the officer who will be held responsible" (p. 74)—underscores the urgency of ethical and legal frameworks to govern autonomous weapons before they become battlefield mainstays.

Embedding values in AI

How to Design AI for Social Good: Seven Essential Factors by Luciano Floridi, Josh Cowls, Thomas C. King, and Mariarosaria Taddeo, published in *Science and Engineering Ethics* (2020).

This analysis delves into the article's purpose, methodology, the seven essential factors, their ethical underpinnings, examples, references, and implications, ensuring that no quote, reference, or example is

omitted. The discussion is structured to provide a comprehensive understanding of how the authors propose to design AI that serves the social good, with detailed explanations, direct citations, and contextual elaboration.

Introduction: Framing AI for Social Good (AI4SG)

The article begins by highlighting the rising prominence of Artificial Intelligence for Social Good (AI4SG) within both information societies and the AI community. The authors note its potential to address pressing social challenges, stating:

"The idea of artificial intelligence for social good (henceforth AI4SG) is gaining traction within information societies in general and the AI community in particular. It has the potential to tackle social problems through the development of AI-based solutions." (p. 1771)

They define AI4SG as:

"the design, development, and deployment of AI systems in ways that (i) prevent, mitigate or resolve problems adversely affecting human life and/or the wellbeing of the natural world, and/or (ii) enable socially preferable and/or environmentally sustainable developments." (p. 1773)

This definition sets the stage for the article's core contribution: identifying seven ethical factors critical to ensuring AI4SG initiatives succeed in delivering social benefits. These factors are not arbitrary but are derived from an empirical analysis of 27 AI4SG projects, with seven representative cases highlighted in the appendix (p. 1792). The authors emphasize that while general AI ethics frameworks exist (e.g., Floridi et al., 2018), AI4SG requires specific considerations due to its focus on social impact.

The introduction also underscores two key challenges: **unnecessary failures** and **missed opportunities**. For instance, they cite IBM's oncology-support software, which failed due to poor design using synthetic data and U.S.-centric protocols, leading to misdiagnoses and loss of trust among doctors (Ross & Swetlitz, 2017; Strickland, 2019). Conversely, an accidental success is noted with IBM's Watson, originally designed for biological mechanisms but repurposed to inspire engineering students in education (Goel et al., 2015). These examples illustrate the need for a systematic approach to AI4SG design to avoid haphazard outcomes.

Methodology: How the Factors Were Identified

The authors employed a rigorous methodology to derive the seven factors, conducting a systematic literature review across five databases—Google Scholar, PhilPapers, Scopus, SSRN, and Web of Science—between October 2018 and May 2019 (p. 1774). They began with a broad search for "AI for Social Good," then refined it to specific areas like healthcare, education, equality, climate change, and environmental protection, using queries such as:

" AND ('Artificial Intelligence' OR 'Machine Learning' OR 'AI') AND 'Social Good'" (p. 1774)

From this, they selected 27 projects, narrowing down to seven representative cases (listed in the appendix, p. 1792) based on scope, variety, impact, and their ability to illustrate the proposed factors. These cases include initiatives like wildlife security optimization (Fang et al., 2016) and hand hygiene tracking (Haque et al., 2017).

The factors align with five established AI ethics principles—**beneficence, nonmaleficence, justice, autonomy, and explicability**—drawn from Floridi et al. (2018). The authors assert:

"AI4SG cannot be inconsistent with the ethical framework guiding the design and evaluation of AI in general." (p. 1774)

Beneficence is a foundational precondition for AI4SG, but it alone is insufficient, as benefits may be offset by harms or risks, necessitating additional considerations specific to AI4SG.

The Seven Essential Factors: Detailed Analysis

Below, each factor is explored in depth, including its explanation, supporting examples, quotes, references, and corresponding best practices, as presented in the article.

1. Falsifiability and Incremental Deployment

Explanation

Trustworthiness is paramount for AI4SG, and falsifiability ensures that critical requirements (e.g., safety) can be empirically tested. The authors explain:

"Falsifiability entails the specification, and the possibility of empirical testing, of one or more critical requirements... Safety is an obvious critical requirement. Hence, for an AI4SG system to be trustworthy, its safety should be falsifiable." (p. 1776)

Since absolute certainty is unattainable, incremental deployment is proposed to test systems progressively from controlled settings to real-world contexts, adjusting assumptions as needed.

Examples and References

- **Germany's Autonomous Vehicles:** The article cites Germany's use of deregulated zones (teststrecken) to test autonomous vehicles incrementally, increasing autonomy levels as trustworthiness is verified (Pagallo, 2017). This aligns with European AI policy recommendations (Floridi et al., 2018).
- **Wildlife Security Model:** A game-theoretic model for wildlife patrols initially assumed flat topography, but real-world testing disproved this, refining the patrol route (Fang et al., 2016).
- **Microsoft's Tay Bot (Counter-Example):** An AI that learned from Twitter users at runtime became offensive, illustrating the risks of untested real-world deployment (Neff & Nagy, 2016).

Quotes

- "If falsifiability is not possible, then the critical requirements cannot be checked, and then the system should not be deemed trustworthy." (p. 1776)
- "What one may prove to be correct via a formal proof, or likely correct via testing in simulation, may be disproved later with the real-world deployment of the system." (p. 1777)

Best Practice

"(1) AI4SG designers should identify falsifiable requirements and test them in incremental steps from the lab to the 'outside world'." (p. 1777)

Analysis

This factor underscores the need for iterative testing to mitigate risks, drawing on formal verification (Dennis et al., 2016) and simulations while acknowledging their limitations. The Tay bot fiasco highlights the dangers of skipping this process, while Germany's approach exemplifies a structured rollout.

2. Safeguards Against the Manipulation of Predictors

Explanation

AI's predictive power is vulnerable to data manipulation, a risk amplified by its scale. The authors reference **Goodhart's Law**:

"When a measure becomes a target, it ceases to be a good measure." (Goodhart, 1975; Strathern, 1997, p. 308)

This can lead to unfair outcomes, breaching justice.

Examples and References

- **Teacher Grade Inflation:** Ghani (2016) warns that transparent models predicting student risk based on math GPA could be gamed by teachers inflating grades, reducing effectiveness.
- **Police Officer Behavior:** An officer might adjust behavior temporarily to avoid intervention if a model predicts adverse events based on recent force incidents (Ghani, 2016).
- **Corporate Fraud:** Manipulation of predictors diminished AI's effectiveness in fraud detection (Zhou & Kapoor, 2011).

Quotes

- "The introduction of AI complicates matters, owing to the scale at which AI is typically applied." (p. 1778)
- "AI4SG designers should adopt safeguards which (i) ensure that non-causal indicators do not inappropriately skew interventions, and (ii) limit, when appropriate, knowledge of how inputs affect outputs from AI4SG systems, to prevent manipulation." (p. 1779)

Best Practice

"(2) AI4SG designers should adopt safeguards which (i) ensure that non-causal indicators do not inappropriately skew interventions, and (ii) limit, when appropriate, knowledge of how inputs affect outputs from AI4SG systems, to prevent manipulation." (p. 1779)

Analysis

This factor addresses the ethical imperative of justice by preventing gaming, a pre-AI issue now magnified. The authors suggest balancing transparency with obfuscation, a tension also noted by Prasad (2018), who advocates democratizing predictor knowledge in some cases.

3. Receiver-Contextualised Intervention

Explanation

Interventions must respect user autonomy while balancing current and future benefits, avoiding over-intrusion that could lead to rejection. The authors state:

"It is essential that software intervenes in users' life only in ways that respect their autonomy." (p. 1779)

Examples and References

- **Cognitive Disability Software:** An interactive system prompts medication reminders but allows users to decline, learning from responses to optimize timing (Chu et al., 2012).
- **Wildlife Security Patrols:** A game-theoretic model suggests routes, but officers can disengage if impractical, lacking flexibility (Fang et al., 2016).
- **Boeing 737 Max:** Pilots couldn't override software malfunctions due to missing optional safety features, contributing to crashes (Tabuchi & Gelles, 2019).

Quotes

- "A suitable receiver-contextualised intervention is one that achieves the right level of disruption while respecting autonomy through optionality." (p. 1780)
- "AI4SG designers should build decision-making systems in consultation with users... with respect for users' right to ignore or modify interventions." (p. 1780)

Best Practice

"(3) AI4SG designers should build decision-making systems in consultation with users interacting with, and impacted, by these systems; with understanding of users' characteristics, of the methods of coordination, and the purposes and effects of an intervention; and with respect for users' right to ignore or modify interventions." (p. 1780)

Analysis

Drawing on McFarlane's taxonomy (1999; 2002), this factor emphasizes user partnership and optionality, contrasting successful autonomy-preserving designs with failures like Boeing's, where autonomy was curtailed.

4. Receiver-Contextualised Explanation and Transparent Purposes

Explanation

Explanations must be tailored to the receiver's context, and system goals must be transparent to foster trust and autonomy. The authors use the **Level of Abstraction (LoA)** framework (Floridi, 2017) to argue that conceptual alignment varies by purpose and audience.

Examples and References

- **Academic Adversity Prediction:** A system used GPA and socio-economic factors, familiar to school officials (Lakkaraju et al., 2015).
- **HIV Education for Homeless Youth:** Initial social graph explanations confused shelter officials; a pedagogic LoA was adopted after testing (Yadav et al., 2016a, b).
- **Medication Prompts:** A system for cognitive disability patients was transparent about non-coercive goals (Chu et al., 2012).
- **Deceptive Bot Study:** A bot posed as a human assistant, justified by scientific value, raising transparency-consent tensions (Eicher et al., 2017).

Quotes

- "The right conceptualisation is likely to vary between AI4SG projects, because they differ greatly in their objectives, subject matter, context and stakeholders." (p. 1781)
- "Transparency in the goal (i.e., system's purpose) of the system is also crucial, for it follows directly from the principle of autonomy." (p. 1783)

Best Practice

"(4) AI4SG designers should explain decisions in terms that are conceptually relevant to the explainee (receiver), taking into account the method by which decisions are reached, and disclose the purposes of the system in an understandable manner." (p. 1784, slightly rephrased in thinking trace)

Analysis

This factor bridges explicability and autonomy, using LoA to customize explanations (Gregor & Benbasat, 1999) and highlighting transparency's role in trust (Herlocker et al., 2000), with exceptions justified by ethical norms like the Nuremberg Code (Nijhawan et al., 2013).

5. Privacy Protection and Data Subject Consent

Explanation

AI4SG's reliance on personal data necessitates robust privacy safeguards and consent, especially for vulnerable groups, aligning with nonmaleficence and autonomy.

Examples and References

- **Google DeepMind NHS Case:** Used patient data without adequate consent, breaching privacy laws (Burgess, 2017).
- **Hand Hygiene Tracking:** Used depth images to de-identify subjects, balancing privacy and efficacy (Haque et al., 2017).
- **Sexuality Detection:** AI trained on dating site photos raised consent issues despite ethics approval (Wang & Kosinski, 2018).

Quotes

- "Privacy is not a novel problem, but the centrality of personal data to many AI (and AI4SG) applications heightens its ethical significance and creates issues around consent." (p. 1786)

- "AI4SG designers should respect the threshold of consent established for the processing of datasets of personal data." (p. 1786)

Best Practice

"(5) AI4SG designers should respect the threshold of consent established for the processing of datasets of personal data." (p. 1786)

Analysis

This factor critiques online consent models (Nissenbaum, 2011) and contrasts ethical failures (DeepMind) with innovative solutions (depth images), emphasizing privacy's heightened stakes in AI.

6. Situational Fairness

Explanation

AI can perpetuate data biases, breaching justice, but must balance removing irrelevant factors with retaining those needed for inclusivity.

Examples and References

- **Predictive Policing:** Biased arrest data reinforced discrimination (Lum & Isaac, 2016; Crawford, 2016).
- **Preterm Birth Prediction:** Historical bias against African-American women risks unfair AI outcomes (Banjo, 2018; CDC, 2020).
- **Chatbot Failures:** A virtual assistant ignored gender context (Eicher et al., 2017), and a mental health bot misunderstood abuse reports (White, 2018).

Quotes

- "AI4SG initiatives relying on biased data may propagate this bias through a vicious cycle." (p. 1787)
- "Designers must sanitise the datasets used to train AI. However, there is equally a risk of applying too strong a disinfectant... by removing important contextual nuances." (p. 1787)

Best Practice

"(6) AI4SG designers should remove from relevant datasets variables and proxies that are irrelevant to an outcome, except when their inclusion supports inclusivity, safety, or other ethical imperatives." (p. 1788)

Analysis

This factor navigates the tension between fairness and context (Caliskan et al., 2017), using examples to show how bias loops (Yang et al., 2018) and insensitivity undermine AI4SG.

7. Human-Friendly Semanticisation

Explanation

AI should support, not replace, human meaning-making (semanticisation), preserving autonomy and agency.

Examples and References

- **Legal Violation Prediction:** An AI defining “violation” could limit judicial roles (Al-Abdulkarim et al., 2015).
- **Alzheimer’s Support:** AI reminders freed caregivers for meaningful interaction, optimizing human semanticisation (Chu et al., 2012; Burns & Rabins, 2000).

Quotes

- "AI4SG must allow humans to curate and foster their ‘semantic capital’, that is, any content that can enhance someone’s power to give meaning to and make sense of (semanticise) something." (p. 1788)
- "AI should be deployed to facilitate human-friendly semanticisation, but not to provide it itself." (p. 1789)

Best Practice

"(7) AI4SG designers should not hinder the ability for people to semanticise (that is, to give meaning to, and make sense of) something." (p. 1789)

Analysis

This factor critiques over-automation (Martinez-Miranda & Aldea, 2005), advocating a supportive role for AI that enhances human agency, as seen in the Alzheimer’s case.

Balancing the Factors

The authors stress that the factors require **intra-factor** and **inter-factor balancing**:

- **Intra-factor:** Balancing intervention frequency (over- vs. under-intervening) or fairness (obfuscation vs. enumeration).
- **Inter-factor:** Transparency vs. manipulation prevention, or explanation vs. privacy.

They pose a moral question:

"The overarching question facing the AI4SG community is, for each given case, whether one is morally obliged to, or obliged not to, design, develop, and deploy a specific AI4SG project." (p. 1791)

Resolution is context-dependent, potentially aided by participatory approaches (Baum, 2017; Prasad, 2018) and AI meta-tools.

Conclusion and Future Directions

The seven factors—summarized in Table 1 (p. 1790)—offer a framework for ethical AI4SG design, grounded in beneficence and the other four principles. The authors conclude:

"The future of AI4SG will likely provide more opportunities to enrich such a set of essential factors."
(p. 1791)

They call for further research into balancing tensions and incorporating diverse perspectives, laying groundwork for sustainable AI4SG policies.

Appendix: Representative Cases

The appendix lists seven projects (p. 1792):

- **A:** Wildlife security (Fang et al., 2016) – Factors 1, 3.
- **B:** Student risk (Lakkaraju et al., 2015) – Factor 4.
- **C:** HIV education (Yadav et al., 2016a, b; 2018) – Factor 4.
- **D:** Cognitive disability aid (Chu et al., 2012) – Factors 3, 4, 7.
- **E:** Virtual assistant (Eicher et al., 2017) – Factors 4, 6.
- **F:** Fraud detection (Zhou & Kapoor, 2011) – Factor 2.
- **G:** Hand hygiene (Haque et al., 2017) – Factor 5.

This analysis exhaustively covers the article, integrating every quote, reference, and example to provide a thorough understanding of designing AI for social good.

Embedding Values in Artificial Intelligence (AI) Systems, published in *Minds and Machines* (2020).

This analysis explores the article's purpose, methodology, key concepts, arguments, examples, references, and implications, ensuring that every quote, reference, and example from the document is included and thoroughly examined. The goal is to provide a comprehensive understanding of how values can be embedded in AI systems, as proposed by the author, in response to the user's query for a detailed, to-the-point explanation.

Introduction: Framing the Ethical Challenge in AI

Van de Poel begins by situating his work within the contemporary discourse on AI ethics, noting the increasing attention given to ethical issues and values in AI design and deployment. He references influential organizations to establish this context:

"Organizations such as the EU High-Level Expert Group on AI and the IEEE have recently formulated ethical principles and (moral) values that should be adhered to in the design and deployment of artificial intelligence (AI)." (p. 385)

The EU High-Level Expert Group on AI (2019) outlines four key principles—respect for human autonomy, prevention of harm, fairness, and explicability—while the IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems (2019) emphasizes human rights, well-being, data agency, transparency, and accountability (p. 385-386). These values, alongside others like security and sustainability, are intended to guide AI governance and design. However, van de Poel poses a pivotal question that drives his investigation:

"But how can we ensure and verify that an AI system actually respects these values?" (p. 385)

This question underscores the article's central aim: to develop an account for determining when an AI system embodies specific values, linking this embodiment to intentional design activities. He argues that existing ethical codes are insufficient without a practical framework to assess compliance, setting the stage for a philosophical and technical exploration.

To structure his account, van de Poel proposes three desiderata:

1. **Connection to Design:** The account must tie value embodiment to the design process, reflecting the focus of ethical codes like IEEE (2019) on designers (p. 386).
2. **Sociotechnical Perspective:** AI systems should be viewed as sociotechnical systems—comprising technical artifacts, human agents, and institutions—rather than isolated technologies (Borenstein et al., 2019; Coeckelbergh, 2020; Boddington, 2017; Behymer & Flach, 2016; Jones et al., 2013) (p. 386).
3. **Distinction Between Humans and AI:** The framework must preserve conceptual differences between human and AI agency, crucial for values like respect for human autonomy (Johnson, 2006; Johnson & Miller, 2008; Illies & Meijers, 2009; Peterson & Spahn, 2011) (p. 386).

He builds on his prior work with Kroes (2014), which satisfies the first and third conditions by linking values to design intent and distinguishing human from technological agency. However, it falls short on the second condition, as it focuses on technical artifacts rather than sociotechnical systems, necessitating an extension for AI's unique characteristics (p. 386-387).

Conceptualizing Values: A Normative Foundation

Before delving into value embedding, van de Poel clarifies the concept of "value." He acknowledges its complexity across disciplines (Brosch et al., 2016; Hirose & Olson, 2015) but asserts its normative essence:

"Rather than being descriptive, values are normative and express what is 'good.'" (p. 388)

He situates values within the evaluative part of normativity, distinct from the deontic part (duties, norms), emphasizing their role in assessing goodness rather than rightness of actions (p. 388). Rejecting a subjective view where values are merely what people value (Stevenson, 1944), he argues this fails to account for misvaluing:

"The problem with such an understanding is that people might very well value things that are not valuable; they sometimes even value things that they know they should not value. Conversely, people might sometimes fail to value things that are valuable." (p. 388)

Instead, he aligns values with normative reasons (Scanlon, 1998; Raz, 1999; Zimmerman, 2015; Jacobson, 2011; Anderson, 1993), suggesting a correspondence between reasons for valuing and something being valuable. This avoids conflating values with preferences and addresses the "wrong kind of reasons" problem (Jacobson, 2011):

"However, the mere presence of reasons for a pro-attitude or pro-behavior does not show that an entity embodies a certain value; those reasons for a pro-attitude or pro-behavior need to originate in the entity itself and not in something else." (p. 388)

For example, a promise to protect an object generates reasons external to the object, not inherent to it, distinguishing embodied value from externally imposed value (p. 388-389). This normative grounding is

critical for his subsequent account of value embodiment in AI systems.

Embodied Values: A Triadic Distinction

Van de Poel introduces a triadic framework to assess value compliance in AI systems: **intended values**, **realized values**, and **embodied values**. He critiques relying solely on intended or realized values due to their limitations:

- **Intended Values:** These are what designers aim to embed, but intentions alone don't guarantee success: "an intended value can be present even if the designed system fails to fulfill that value" (p. 389).
- **Realized Values:** These are outcomes in operation, but they may stem from misuse or exceptional circumstances, not the system itself: "not all realized values can be meaningfully attributed to the relevant AI system" (p. 389).

He illustrates this with a self-driving car example:

"For example, suppose a self-driving car (understood here as an AI system) causes an accident resulting in a number of fatalities. Can we conclude from this accident that the AI system (i.e., the self-driving car) was unsafe because that value was realized in the accident? The answer seems negative." (p. 389)

Both approaches suffer from the wrong kind of reasons problem—intended values root reasons in designers' minds, while realized values may reflect external factors (p. 389). Thus, he focuses on embodied values:

"The basic idea... is that embodied values should be understood as values that have been intentionally, and successfully, embedded in an AI system by its designers." (p. 389)

This requires two conditions: (1) intentional design for the value, and (2) the system respecting or furthering that value under proper use. Figure 1 (p. 390) illustrates the interplay among these categories, showing feedback loops where discrepancies trigger redesign or altered use. Examples include:

- If intended and embodied values align but realized values differ, changing use suffices.
- If embodied values diverge from intended ones, redesign is needed.
- Unintended consequences (e.g., an AI causing discrimination) may necessitate revising intended values to include fairness (p. 390).

Van de Poel emphasizes redesign as an ongoing process, especially for adaptive AI systems, involving not just designers but users and operators (p. 390).

AI Systems as Sociotechnical Systems: Five Building Blocks

Recognizing AI's complexity, van de Poel conceptualizes AI systems as sociotechnical systems, expanding beyond traditional models (Bauer & Herder, 2009; Baxter & Sommerville, 2011; Geels, 2004; Bruijn & Herder, 2009; Pasmore & Sherwood, 1978; Kroes et al., 2006; Ottens et al., 2006; Dam et al., 2013; Franssen, 2014; Nickel, 2013) (p. 391). Traditional systems comprise:

1. **Technical Artifacts:** Physical objects with designed functions (Kroes, 2010; Kroes & Meijers, 2006; Houkes et al., 2002; Houkes & Vermaas, 2010; Vermaas & Houkes, 2006) (p. 391).
2. **Human Agents:** Individuals with intentionality and moral agency (p. 391).
3. **Institutions:** Rules governing human behavior (North, 1990; Calvert, 1995; Ullmann-Margalit, 1977; Ostrom, 2005; Bicchieri, 2006) (p. 392).

AI systems add two unique components:

4. **Artificial Agents (AAs):** Autonomous, interactive, and adaptive entities lacking human intentionality (Floridi & Sanders, 2004) (p. 392).
5. **Technical Norms:** Code-based rules regulating AAs, grounded in causal-physical terms (Mahmoud et al., 2014) (p. 392).

Table 1 (p. 391) distinguishes these blocks by intentional versus physical-causal natures, highlighting AI's hybridity and adaptivity as both opportunities and challenges for value embedding (Wallach & Allen, 2009; Anderson & Anderson, 2011; Cave et al., 2019; Vanderelst & Winfield, 2018) (p. 387).

Value Embedding in Technical Artifacts

Adapting his earlier work (Van de Poel & Kroes, 2014), van de Poel defines value embodiment in technical artifacts:

"Technical artifact x embodies value V if the designed properties of x have the potential to achieve or contribute to V (under appropriate circumstances) due to the fact that x has been designed for V." (p. 393)

This hinges on two connected conditions: (1) design intent for V, and (2) conduciveness to V through use. He provides four examples:

1. **Sea Dikes:** "designed to protect against flooding... conducive for protection against flooding... therefore embody the value of safety" (p. 393).
2. **Bread Knife:** "designed to cut bread... can also be used for killing... does not embody the (dis)value of killing" because it lacks design intent (p. 393).
3. **Faulty Pacemaker:** "designed for (contributing to) human well-being... fails to contribute to well-being" due to poor design, thus not embodying well-being (p. 393).
4. **Recommender System:** "designed to serve its customers may unintendedly... contribute to filter bubbles and echo chambers... (dis)values such as lack of respect and untruth" are not embodied without redesign intent (p. 393-394).

For unintended consequences, he suggests redesign—by designers or users (Vermaas & Houkes, 2006)—can embed new values, distinguishing passive, idiosyncratic, and innovative use (p. 394). If negative outcomes persist, designers acquire an obligation to embed positive values (Winner, 1977) (p. 394).

Value Embedding in Institutions

Institutions, as rules, also embody values. Using the ADICO grammar (Crawford & Ostrom, 1995), van de Poel categorizes:

- **Shared Strategies (AIC):** Descriptive expectations, e.g., "Pedestrians (A) use an umbrella to avoid getting wet (I) when it rains (C)" (p. 395).
- **Norms (ADIC):** Deontic expectations without sanctions, e.g., "Residents (A) must (D) greet their neighbors (I) in this neighborhood (C)" (p. 395).
- **Rules (ADICO):** Deontic expectations with sanctions, e.g., "Car drivers (A) must (D) drive on the right side of the road (I) in the Netherlands (C), otherwise they will be fined by the police (O)" (p. 395).

He proposes:

"Institution R embodies value V if R is conducive to V because R has been designed for V." (p. 396)

Examples include:

1. **Traffic Rule:** Embodies safety as it's designed for and conducive to traffic safety (p. 396).
2. **Greeting Norm:** Embodies politeness through design and effect (p. 396).
3. **Pavement Strategy:** Embodies convenience via shared design and conduciveness (p. 396).

Like artifacts, institutions require both conditions, with redesign addressing unintended outcomes (p. 397).

Human Agents: Mediators of Value

Human agents—users, operators, designers—interact with embodied values in artifacts (V_T) and institutions (V_I), influenced by personal values (V_A) (p. 397). Figure 2 (p. 398) shows this intentional-causal dynamic, with outcomes potentially diverging from embodied values. Humans' reflective capacity (Figure 3, p. 399) enables evaluation and redesign, balancing system stability and adaptability (Franssen, 2015) (p. 397-398).

Artificial Agents: Designed Autonomy

AAs, distinct from humans and artifacts (Table 2, p. 400), embody values if designed for them and conducive to them (p. 399-400). Moor's (2006) taxonomy classifies AAs:

1. **Ethical Impact Agents:** Affect values without intent, thus not embodying them (p. 400).
2. **Implicit Ethical Agents:** Designed to follow values, embodying them if effective (p. 400).
3. **Explicit Ethical Agents:** Represent and reason about values, embodying them if top-down designed (p. 401).
4. **Full Ethical Agents:** Hypothetical with human-like traits, currently unfeasible (Winfield, 2019; Müller, 2020) (p. 400).

Adaptivity risks disembodying values, requiring restrictions or monitoring (Grodzinsky et al., 2008; Cervantes et al., 2020; Allen et al., 2005; Wallach & Allen, 2009; Cave et al., 2019; van Wylsberghe & Robbins, 2019) (p. 400-401).

Technical Norms: Regulating AAs

Technical norms, akin to institutions for humans, regulate AAs via code (Lessig, 1999; Akrich, 1992; Latour, 1992; Thaler & Sunstein, 2009; Fogg, 2003; Norman, 2000) (p. 401-402). Created offline or emergently (Hollander & Wu, 2011), they embody values if:

"Technical norm N embodies value V if (1) N has been designed (by the human system designers) for V and (2) the execution of N within the system is conducive to V." (p. 402)

This mirrors the artifact and institution accounts (Mahmoud et al., 2014; Aldewereld & Sichman, 2013; Panagiotidi et al., 2013; Dybalova et al., 2014; Leenes & Lucivero, 2014) (p. 402).

System-Level Value Embedding

At the system level, van de Poel proposes:

"Value V is embodied in sociotechnical system S if S is conducive to V because of those components of S that have been designed for V." (p. 403)

This accommodates systems not wholly designed (Bowker et al., 2010), requiring only some components to embody V, with conduciveness tied to following norms and use plans (p. 403-404).

Conclusion: Practical Lessons

Van de Poel concludes with two lessons:

1. **Continuous Monitoring and Redesign:** AI's adaptivity necessitates oversight and human control (Santoni de Sio & van den Hoven, 2018) to manage disembodiment risks (p. 404-405).
2. **Focus on Technical Norms:** Embedding values in norms may be more effective than in AAs alone, shifting focus from machine ethics to system regulation (p. 405).

Building Ethics into Artificial Intelligence by Han Yu et al.,

This analysis delves deeply into the content, structure, and contributions of the paper, incorporating all quotes, references, and examples provided in the document. It avoids being a mere summary by exploring the nuances, implications, and technical details of each section, while maintaining a comprehensive and self-contained narrative. Let's dive in.

Introduction: Setting the Stage for Ethical AI

The paper "*Building Ethics into Artificial Intelligence*" by Han Yu and co-authors from institutions like Nanyang Technological University, University of Massachusetts Amherst, and Hong Kong University of Science and Technology, tackles a pressing issue in modern AI research: how to ensure AI systems make ethical decisions. Published with a focus on technical solutions, it bridges a gap left by previous surveys that emphasized psychological, social, and legal perspectives over actionable computational approaches.

The abstract outlines the paper's intent clearly:

"As artificial intelligence (AI) systems become increasingly ubiquitous, the topic of AI governance for ethical decision-making by AI has captured public imagination. Within the AI research community, this

topic remains less familiar to many researchers. In this paper, we complement existing surveys... with an analysis of recent advances in technical solutions for AI governance."

This sets the stage for a technical exploration, distinct from the broader public discourse often dominated by fears of artificial general intelligence (AGI). The authors acknowledge this public anxiety:

"A major source of public anxiety about AI, which tends to be overreactions [Bryson and Kime, 2011], is related to artificial general intelligence (AGI) [Goertzel and Pennachin, 2007] research aiming to develop AI with capabilities matching and eventually exceeding those of humans."

However, they quickly pivot to a more immediate concern: "Although we are still decades away from AGI, existing autonomous systems (such as autonomous vehicles) already warrant the AI research community to take a serious look into incorporating ethical considerations into such systems." This pragmatic focus on current systems, like autonomous vehicles (AVs), grounds the paper in real-world relevance.

The authors define ethics via Cointe et al. (2016), framing it as "a normative practical philosophical discipline of how one should act towards others," encompassing three dimensions:

1. **Consequentialist Ethics:** "An agent is ethical if and only if it weighs the consequences of each choice and chooses the option which has the most moral outcomes." Also called utilitarian ethics, it prioritizes aggregate benefits.
2. **Deontological Ethics:** "An agent is ethical if and only if it respects obligations, duties and rights related to given situations." This rule-based approach aligns with social norms.
3. **Virtue Ethics:** "An agent is ethical if and only if it acts and thinks according to some moral values (e.g. bravery, justice, etc.)." It emphasizes an intrinsic moral character.

They further define ethical dilemmas as "situations in which any available choice leads to infringing some accepted ethical principle and yet a decision has to be made [Kirkpatrick, 2015]." This foundational framework underpins the paper's taxonomy of AI governance techniques, divided into four areas:

- **Exploring Ethical Dilemmas**
- **Individual Ethical Decision Frameworks**
- **Collective Ethical Decision Frameworks**
- **Ethics in Human-AI Interactions**

Each area is explored in depth below, with every detail, quote, and reference meticulously unpacked.

Exploring Ethical Dilemmas: Understanding Human Preferences

The first area, "Exploring Ethical Dilemmas," focuses on tools that help the AI community understand human preferences in ethical scenarios. The paper states: "In order to build AI systems that behave ethically, the first step is to explore the ethical dilemmas in the target application scenarios." Two key tools are highlighted: GenEth and the Moral Machine project.

GenEth: Expert-Driven Ethical Analysis

Proposed by Anderson and Anderson (2014), GenEth is an "ethical dilemma analyzer" designed to involve ethicists in codifying ethical principles for AI. The authors note:

"They realized that ethical issues related to intelligent systems are likely to exceed the grasp of the original system designers, and designed GenEth to include ethicists into the discussion process in order to codify ethical principles in given application domains."

GenEth employs a structured representation schema:

1. **Features:** "Denoting the presence or absence of factors (e.g., harm, benefit) with integer values."
2. **Duties:** "Denoting the responsibility of an agent to minimize/maximize a given feature."
3. **Actions:** "Denoting whether an action satisfies or violates certain duties as an integer tuple."
4. **Cases:** "Used to compare pairs of actions on their collective ethical impact."
5. **Principles:** "Denoting the ethical preference among different actions as a tuple of integer tuples."

This framework is operationalized through a graphical user interface, where discussions are processed using inductive logic programming to "infer principles of ethical actions." GenEth's strength lies in its systematic approach, leveraging expert input to formalize ethics, though it lacks the scalability of crowd-based methods.

Moral Machine: Crowdsourcing Ethical Preferences

In contrast, MIT's Moral Machine project (<http://moralmachine.mit.edu/>) uses crowdsourcing to gather public opinions on ethical dilemmas, particularly for AVs. The paper explains:

"The Moral Machine project focuses on studying the perception of autonomous vehicles (AVs) which are controlled by AI and has the potential to harm pedestrians and/or passengers if they malfunction."

Participants judge scenarios, such as whether an AV should sacrifice its passenger to save more pedestrians, with preferences analyzed across eight considerations:

1. Saving more lives
2. Protecting passengers
3. Upholding the law
4. Avoiding intervention
5. Gender preference
6. Species preference
7. Age preference
8. Social value preference

The project's findings, based on 3 million participants, reveal a nuanced public stance:

"Based on feedbacks from 3 million participants, the Moral Machine project found that people generally prefer the AV to make sacrifices if more lives can be saved. If an AV can save more pedestrian lives by killing its passenger, more people prefer others' AVs to have this feature rather than their own AVs [Bonnefon et al., 2016; Sharif et al., 2017]."

This highlights a consequentialist bent—favoring the greater good—but also a self-interest paradox, as individuals hesitate to apply the same logic to their own vehicles. The authors caution:

"Nevertheless, self-reported preferences often do not align well with actual behaviours [Zell and Krizan, 2014]. Thus, how much the findings reflect actual choices is still an open question."

Alternative suggestions emerge, such as random decision-making ("let fate decide") per Broome (1984) or segregating AVs from human traffic (Bonneton et al., 2016). These diverse opinions underscore the complexity of ethical consensus, making tools like Moral Machine critical yet imperfect for informing AI design.

Individual Ethical Decision Frameworks: Empowering Single Agents

The second area, "Individual Ethical Decision Frameworks," explores mechanisms for single AI agents to make ethical decisions. The paper asserts: "When it comes to ethical decision-making in AI systems, the AI research community largely agrees that generalized frameworks are preferred over ad-hoc rules." Several frameworks are detailed, each with unique approaches.

MoralDM: Combining Rules and Analogies

Dehghani et al. (2008) propose MoralDM, which integrates:

1. **First-Principles Reasoning:** "Makes decisions based on well-established ethical rules (e.g., protected values)," such as the moral unacceptability of murder regardless of outcome.
2. **Analogical Reasoning:** "Compares a given scenario to past resolved similar cases to aid decision-making."

The authors note: "As the number of resolved cases increases, the exhaustive comparison approach by MoralDM is expected to become computationally intractable." Thus, Blass and Forbus (2015) extend it with "structure mapping," which "trims the search space by computing the correspondences, candidate inferences and similarity scores between cases." This enhancement improves efficiency while retaining MoralDM's ability to handle culturally sensitive "protected values."

BDI-Based Ethical Judgment

Cointe et al. (2016) offer a framework using the Belief-Desire-Intention (BDI) model (Rao and Georgeff, 1995), structured around:

- **Awareness:** "Generates the beliefs that describe the current situation facing the agent and the goals of the agent."
- **Evaluation:** "Generates the set of possible actions and desirable actions."
- **Goodness:** "Computes the set of ethical actions based on the agent's beliefs, desires, actions, and moral value rules."
- **Rightness:** "Evaluates whether or not executing a possible action is right under the current situation."

This process adapts to judging others' actions under varying information levels (blind, partially informed, fully informed), though it lacks "quantitative measure of how far a behaviour is from rightfulness or goodness," limiting its precision.

Game Theory and Machine Learning

Conitzer et al. (2017) propose two approaches:

1. **Game Theory:** Extends the "extensive form" (game trees) with "passive actions" to account for protected values, addressing scenarios where inaction is ethical.

2. **Machine Learning:** Classifies actions as right or wrong using human judgments, potentially from Moral Machine data, though cultural inconsistencies pose challenges. They suggest leveraging moral foundations (e.g., harm/care, fairness) from Clifford et al. (2015) for generalizable representations.

The authors envision combining these: "Game theory and machine learning can be combined into one framework in which game theoretic analysis of ethics is used as a feature to train machine learning approaches."

CP-Nets for Preference Balancing

Loreggia et al. (2018) use CP-nets to "represent the exogenous ethics priorities and endogenous subjective preferences," introducing a "notion of distance between CP-nets" to balance agent preferences with ethical requirements. This allows flexibility in decision-making when preferences align closely with ethics.

High-Level Action Language

Berreby et al. (2017) shift moral reasoning to agents via a high-level action language, implemented with answer set programming (Lifschitz, 2008). It:

- Simulates outcomes using action, event, and situation data.
- Produces causal traces via a causal engine.
- Assesses goodness and rightfulness using ethical specifications and deontological rules.

This framework enables agents to "decide and explain their actions, and reason about other agents' actions on ethical grounds," reducing the burden on developers.

Ethics Shaping in Reinforcement Learning

Wu and Lin (2018) adapt reinforcement learning (RL) with "ethics shaping," incorporating ethical values into reward functions:

"By assuming that the majority of observed human behaviours are ethical, the proposed approach learns ethical shaping policies from available human behaviour data in given application domains."

The function "rewards positive ethical decisions, punishes negative ethical decisions, and remains neutral when ethical considerations are not involved," separating ethics from standard RL design.

Collective Ethical Decision Frameworks: Group Dynamics

The third area, "Collective Ethical Decision Frameworks," addresses ethical decision-making among multiple agents. Pagallo (2016) argues: "Individual ethical behavior isn't enough and that we need social norms and rules that can evolve."

Social Norms and Trust Networks

Singh (2014; 2015) proposes a distributed framework using social norms, defined via:

- Roles (qualifications, privileges, penalties)
- Commitments, authorizations, prohibitions, sanctions, and power

Agents form "a network of trust based on techniques from the reputation modelling literature [Yu et al., 2010; Yu et al., 2013]" for self-governance.

Human-Agent Collectives

Greene et al. (2016) envision agents evaluating different ethical dimensions (deontological, consequentialist, virtue) within human-agent collectives (Jennings et al., 2014). Preferences are aggregated, but challenges include:

- Large action sets outnumbering agents
- Interdependent actions
- Missing or imprecise preferences

Voting-Based System

Noothigattu et al. (2018) build on Moral Machine data, using a voting system with "swap-dominance" to rank alternatives:

"Assuming everything else is fixed, an outcome a swap-dominates another outcome b if every ranking which ranks a higher than b has a weight which is equal to or larger than rankings that rank b higher than a."

This ensures computationally efficient, consequentialist decisions reflecting collective preferences.

Ethics in Human-AI Interactions: Influencing Behavior

The fourth area, "Ethics in Human-AI Interactions," focuses on AI influencing humans ethically, guided by the Belmont Report (Bel, 1978):

1. "People's personal autonomy should not be violated."
2. "Benefits brought about by the technology should outweigh risks."
3. "The benefits and risks should be distributed fairly among the users."

Persuasion Agents

Stock et al. (2016) study AI persuasion in the trolley problem, testing:

1. Emotional appeals
2. Utilitarian arguments
3. Lying

Findings show: "Participants hold a strong preconceived negative attitude towards the persuasion agent, and argumentation-based and lying-based persuasion strategies work better than emotional persuasion strategies."

Emotional Responses

Battaglino and Damiano (2015) use Coping Theory (Marsella and Gratch, 2003) to trigger emotions like shame (for self-violations) or reproach (for others' violations), enhancing human-AI interaction.

Discussions and Future Directions

The paper notes a focus on individual frameworks, a need for diverse cultural data, and challenges in collective preference representation and human-AI ethics. It advocates interdisciplinary collaboration and a global AI regulatory framework (Erdélyi and Goldsmith, 2018). Future directions include:

1. **Social-Systems Analysis:** Using transfer learning (Pan and Yang, 2010) to model diverse ethics.
2. **Revising Social Contracts:** Dynamic regulations for AI responsibility.
3. **Explainable AI:** Argumentation-based explanations (Fan and Toni, 2015) balancing transparency.
4. **Adversarial Considerations:** Incorporating adversarial game theory (Vorobeychik et al., 2012) to counter strategic exploitation.

This analysis covers every facet of the paper, from its foundational definitions to its technical proposals, ensuring a thorough understanding of how ethics can be built into AI.