

**“I, Robot – I, Criminal”—When Science Fiction Becomes Reality:
Legal Liability of AI Robots committing Criminal Offenses**

Gabriel Hallevy*

I. INTRODUCTION

Can society impose criminal liability upon robots? The technological world has changed rapidly. Simple human activities are being replaced by robots. As long as humanity used robots as mere tools, there was no real difference between robots and screwdrivers, cars or telephones. When robots became sophisticated, we used to say that robots “think” for us. The problem began when robots evolved from “thinking” machines into thinking machines (without quotation marks)—or Artificial Intelligence Robots (AI Robots). Could they become dangerous?

Unfortunately, they already are. In 1950, Isaac Asimov set down three fundamental laws of robotics in his science fiction masterpiece “I, Robot”: (1) a robot may not injure a human being or, through inaction, allow a human being to come to harm; (2) a robot must obey the orders given to it by human beings, except where such orders would conflict with the First Law; and (3) a robot must protect its own existence, as long as such protection does not conflict with the First or Second Laws.¹ These three fundamental laws are obviously contradictory.² What if a man orders a robot to hurt another person for the own good of the other person? What if the

* Associate Professor, Faculty of Law, Ono Academic College.

¹ ISSAC ASIMOV, I, ROBOT (1950).

² Isaac Asimov wrote in his introduction to THE REST OF ROBOTS (1964) that “[t]here was just enough ambiguity in the Three Laws to provide the conflicts and uncertainties required for new stories, and, to my great relief, it seemed always to be possible to think up a new angle out of the 61 words of the Three Laws.”

robot is in police service and the commander of the mission orders it to arrest a suspect and the suspect resists arrest? Or what if the robot is in medical service and is ordered to perform a surgical procedure on a patient, the patient objects, but the medical doctor insists that the procedure is for the patient's own good, and repeats the order to the robot?

The main question in that context is which kind of laws or ethics are correct and who is to decide. In order to cope with these same problems as they relate to humans, society devised criminal law. Criminal law embodies the most powerful legal social control in modern civilization. People's fear of AI robots, in most cases, is based on the fact that AI robots are not considered to be subject to the law, specifically to criminal law. In the past, people were similarly fearful of corporations and their power to commit a spectrum of crimes, but since corporations are legal entities subject to criminal and corporate law, that kind of fear has been reduced significantly.³

The apprehension that AI robots evoke may have arisen due to Hollywood's depiction of AI robots in numerous films, such as "2001: A Space Odyssey,"⁴ and the modern trilogy "The Matrix,"⁵ in which AI robots are not subject to the law. However, it should be noted that Hollywood did treat AI robots in an empathic way by depicting them as human, as almost

³ See generally John C. Coffee, Jr., "No Soul to Damn: No Body to Kick": An Unscandalised Inquiry Into the Problem of Corporate Punishment, 79 MICH. L. REV. 386 (1981); STEVEN BOX, POWER, CRIME AND MYSTIFICATION 16-79 (1983); Brent Fisse & John Braithwaite, *The Allocation of Responsibility for Corporate Crime: Individualism, Collectivism and Accountability*, 11 SYDNEY L. REV. 468 (1988).

⁴ STANLEY KUBRICK, 2001: A SPACE ODYSSEY (1968).

⁵ JOEL SILVER, THE MATRIX (1999); JOEL SILVER, LAURENCE WACHOWSKI AND ANDREW PAUL WACHOWSKI, THE MATRIX RELOADED (2003); JOEL SILVER, LAURENCE WACHOWSKI AND ANDREW PAUL WACHOWSKI, THE MATRIX REVOLUTIONS (2003).

human, or as wishing to be human.⁶ This kind of treatment included, of course, clear subordination to human legal social control and to criminal law.

The modern question relating to AI robots becomes: Does the growing intelligence of AI robots subject them to legal social control, just as any other legal entity?⁷ This article attempts to work out a legal solution to the problem of the criminal liability of AI robots. At the outset, a definition of an AI robot will be presented. Based on that definition, this article will then propose and introduce three models of AI robot criminal liability: (1) the perpetration-by-another liability model, (2) the natural-probable-consequence liability model, and (3) the direct liability model.

These three models might be applied separately, but in many situations, a coordinated combination of them (all or some of them) is required in order to complete the legal structure of criminal liability. Once we examine the possibility of legally imposing criminal liability on AI robots, then the question of punishment must be addressed. How can an AI robot serve a sentence of imprisonment? How can capital punishment be imposed on an AI robot? How can probation, a pecuniary fine, or the like be imposed on an AI robot? Consequently, it is necessary to formulate viable forms of punishment in order to impose criminal liability practically on AI robots.

⁶ See, e.g., STEVEN SPIELBERG, STANLEY KUBRICK, JAN HARLAN, KATHLEEN KENNEDY, WALTER F. PARKES AND BONNIE CURTIS, *A.I. ARTIFICIAL INTELLIGENCE* (2001).

⁷ See in general, but not in relation to criminal law, e.g., Thorne L. McCarty, *Reflections on Taxman: An Experiment in Artificial Intelligence and Legal Reasoning*, 90 HARV. L. REV. 837 (1977); Donald E. Elliott, *Holmes and Evolution: Legal Process as Artificial Intelligence*, 13 J. LEGAL STUD. 113 (1984); Thomas E. Headrick & Bruce G. Buchanan, *Some Speculation about Artificial Intelligence and Legal Reasoning*, 23 STAN. L. REV. 40 (1971); Antonio A. Martino, *Artificial Intelligence and Law*, 2 INT'L J.L. & INFO. TECH. 154 (1994); Edwina L. Rissland, *Artificial Intelligence and Law: Stepping Stones to a Model of Legal Reasoning*, 99 YALE L.J. 1957 (1990).

II. WHAT IS AN AI ROBOT?

For some years, there has been significant controversy about the very essence of AI robots.⁸ Futurologists have proclaimed the birth of a new species, *machina sapiens*, which will share the human place as intelligent creatures on Earth.⁹ Critics have argued that a “thinking machine” is an oxymoron.¹⁰ Machines, including robots, with their foundations of cold logic, can never be insightful or creative as humans are.¹¹ This controversy raises the basic questions of the essence of humanity (*i.e.*, do human beings function as thinking machines?) and of AI (*i.e.*, can there be thinking machines?).¹²

There are five attributes that one would expect an intelligent entity to have:¹³

(i) communication (One can communicate with an intelligent entity. The easier it is to communicate with an entity, the more intelligent the entity seems. One can communicate with a dog, but not about Einstein’s theory of relativity. One can communicate with a little child about Einstein’s theory, but it requires a discussion in terms that a child can comprehend.); **(ii) mental knowledge** (An intelligent entity is expected to have some knowledge about itself.); **(iii)**

⁸ Terry Winograd, *Thinking Machines: Can There Be? Are We?*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 167 (Derek Partridge & Yorick Wilks eds., 2006).

⁹ *Id.*

¹⁰ *Id.*

¹¹ *Id.*

¹² For the formal foundations of AI, *see generally* Teodor C. Przymusinski, *Non-Monotonic Reasoning versus Logic Programming: A New Perspective*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 49 (Derek Partridge & Yorick Wilks eds., 2006); Richard W. Weyhrauch, *Prolegomena to a Theory of Mechanized formal Reasoning*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 72 (Derek Partridge & Yorick Wilks eds., 2006).

¹³ Roger C. Schank, *What is AI, Anyway?*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 3 (Derek Partridge & Yorick Wilks eds., 2006).

external knowledge (An intelligent entity is expected to know about the outside world, to learn about it, and utilize that information.); **(iv) goal-driven behavior** (An intelligent entity is expected to take action in order to achieve its goals.); and **(v) creativity** (An intelligent entity is expected to have some degree of creativity. In this context, creativity means the ability to take alternate action when the initial action fails. A fly that tries to exit a room and bumps into a window pane, tries to do that over and over again. When an AI robot bumps into a window, it tries to exit using the door).

Most AI robots possess these five attributes by definition.¹⁴ Some 21st Century types of AI robots possess even more attributes that enable them to act in far more sophisticated ways. In November 2009, during the Supercomputing Conference in Portland Oregon (“SC 09”), IBM scientists and others announced that they succeeded in creating a new algorithm named “Blue Matter,” which possesses the thinking capabilities of a cat.¹⁵ This algorithm collects information from many units with parallel and distributed connections.¹⁶ The information is integrated and creates a full image of sensory information, perception, dynamic action and reaction, and cognition.¹⁷ This platform simulates brain capabilities, and eventually, it is supposed to simulate

¹⁴ Schank, *supra* note 13, at 4-6.

¹⁵ Chris Capps, “Thinking” Supercomputer Now Conscious as a Cat, UNEXPLAINABLE.NET, Nov. 19, 2009, http://www.unexplainable.net/artman/publish/article_14423.shtml; *see also* Super Computing, <http://sc09.supercomputing.org>.

¹⁶ *Id.*

¹⁷ *Id.*

real thought processes.¹⁸ The final application of this algorithm contains not only analog and digital circuits, metal or plastics, but also protein-based biologic surfaces.¹⁹

An AI robot has a wide variety of applications.²⁰ A robot can be designed to imitate the physical capabilities of a human being, and these capabilities can be improved.²¹ For instances, a robot is capable of being physically faster and stronger than a human being.²² The AI software installed in it also enables the robot to calculate many complicated calculations faster and simultaneously, or to “think” faster.²³ An AI robot is capable of learning and of gaining experience, and experience is a useful way of learning.²⁴ All these attributes create the essence of an AI robot. AI robots and AI software are used in a wide range of applications in industry, military services, medical services, science, and even in games.²⁵

¹⁸ Capps, *supra* note 15.

¹⁹ *Id.*

²⁰ See, e.g., Yorick Wilks, *One Small Head: Models and Theories*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 121 (Derek Partridge & Yorick Wilks eds., 2006); Alan Bundy & Stellan Ohlsson, *The Nature of AI Principles*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 135 (Derek Partridge & Yorick Wilks eds., 2006); Thomas W. Simon, *Artificial Methodology Meets Philosophy*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 155 (Derek Partridge & Yorick Wilks eds., 2006).

²¹ *Id.*

²² *Id.*

²³ *Id.*

²⁴ *Id.*

²⁵ See, e.g., William B. Schwartz, Ramesh S. Patil & Peter Szolovits, *Artificial Intelligence in Medicine Where Do We Stand*, 27 JURIMETRICS J. 362 (1987); Richard E. Susskind, *Artificial Intelligence, Expert Systems and the Law*, 5 DENNING L.J. 105 (1990).

III. MODELS OF THE CRIMINAL LIABILITY OF AI ROBOTS

The fundamental question of criminal law is the question of criminal liability, *i.e.*, whether the specific entity (human or corporation) bears criminal liability for a specific offense committed at a specific point in time and space.²⁶ In order to impose criminal liability upon a person, two main elements must exist.²⁷ The first is the factual element, *i.e.*, criminal conduct (*actus reus*), while the other is the mental element, *i.e.*, knowledge or general intent in relation to the conduct element (*mens rea*).²⁸ If one of them is missing, no criminal liability can be imposed.²⁹ The *actus reus* requirement is expressed mainly by acts or omissions.³⁰ Sometimes, other factual elements are required in addition to conduct, such as the specific results of that conduct and the specific circumstances underlying the conduct.³¹ The *mens rea* requirement has various levels of mental elements.³² The highest level is expressed by knowledge, while

²⁶ See generally JEROME HALL, GENERAL PRINCIPLES OF CRIMINAL LAW 70-211 (2d ed. 2005) (1960).

²⁷ *Id.*

²⁸ *Id.*

²⁹ *Id.*

³⁰ See Walter Harrison Hitchler, *The Physical Element of Crime*, 39 DICK. L. REV. 95 (1934); MICHAEL MOORE, ACT AND CRIME: THE PHILOSOPHY OF ACTION AND ITS IMPLICATIONS FOR CRIMINAL LAW 5 (1993) (Tony Honore & Joseph Raz, eds., Oxford University Press 1993).

³¹ See JOHN WILLIAM SALMOND, ON JURISPRUDENCE 505 (Glanville Williams ed., 11th ed. 1957); GLANVILLE WILLIAMS, CRIMINAL LAW: THE GENERAL PART 18 (2d ed., Steven & Sons Ltd. 1961); OLIVER W. HOLMES, THE COMMON LAW 54 (Mark DeWolf Howe ed., Harvard University Press 1923) (1881); Walter Wheeler Cook, *Act, Intention, and Motive in Criminal Law*, 26 YALE L.J. 645 (1917).

³² HALL, *supra* note 26, at 105-45, 325-59.

sometimes it is accompanied by a requirement of intent or specific intention.³³ Lower levels are expressed by negligence (a reasonable person should have known),³⁴ or by strict liability offenses.³⁵

No other criteria or capabilities are required in order to impose criminal liability, not from humans, nor from any other kind of entity, including corporations and AI robots.³⁶ An entity might possess further capabilities; however, in order to impose criminal liability, the existence of *actus reus* and *mens rea* in the specific offense is quite enough.³⁷ As far as known to science, a spider is capable of acting, but it is incapable of formulating the *mens rea* requirement; therefore, a spider bite bears no criminal liability. A parrot is capable of repeating words it hears, but it is incapable of formulating the *mens rea* requirement for libel. In order to impose criminal liability on any kind of entity, it must be proven that the above two elements

³³ See generally J. Ll. J. Edwards, *The Criminal Degrees of Knowledge*, 17 MOD. L. REV. 294, 295 (1954); Rollin M. Perkins, “*Knowledge*” as a *Mens rea* Requirement, 29 HASTINGS L.J. 953 (1978); and see, e.g., *United States v. Youts*, 229 F.3d 1312, 1316 (10th Cir. 2000); *United States v. Spinney*, 65 F.3d 231, 235 (1st Cir. 1995); *People v. Steinberg*, 595 N.E.2d 845, 847 (N.Y. 1992); *State v. Sargent*, 594 A.2d 401, 403 (Vt. 1991); *State v. Wyatt*, 482 S.E.2d 147, 150 (W. Va. 1996).

³⁴ See, e.g., Jerome Hall, *Negligent Behaviour Should Be Excluded from Penal Liability*, 63 COLUM. L. REV. 632, 632 (1963); Robert P. Fine & Gary M. Cohen, *Is Criminal Negligence a Defensible Basis for Criminal Liability?*, 16 BUFF. L. REV. 749, 780 (1966).

³⁵ See, e.g., Jeremy Horder, *Strict Liability, Statutory Construction and the Spirit of Liberty*, 118 L. Q. REV. 458, 458 (2002); Francis Bowes Sayre, *Public Welfare Offenses*, 33 COLUM. L. REV. 55, 55 (1933); Stuart P. Green, *Six Senses of Strict Liability: A Plea for Formalism*, in APPRAISING STRICT LIABILITY 1 (A. P. Simester ed., 2005); A. P. Simester, *Is Strict Liability Always Wrong?*, in APPRAISING STRICT LIABILITY 21 (A. P. Simester ed., 2005).

³⁶ HALL, *supra* note 26.

³⁷ *Id.* at 185-86.

existed.³⁸ Thus, when it has been proven that a person committed the criminal act knowingly or with criminal intent, that person is held criminally liable for that offense.³⁹ The relevant question concerning the criminal liability of AI robots is: How can these entities fulfill the two requirements of criminal liability? This article proposes the imposition of criminal liability on AI robots using three possible models of liability: (A) the Perpetration-by-Another liability model; (B) the Natural-Probable-Consequence liability model; and (C) the Direct liability model. The following is an explanation of these possible models:

A. The Perpetration-by-Another Liability Model: AI Robots as Innocent Agents

This first model does not consider the AI robot as possessing any human attributes. The AI robot is considered an innocent agent. Accordingly, due to that legal viewpoint, a machine is a machine, and is never human. However, one cannot ignore an AI robot's capabilities, as previously mentioned.⁴⁰ Pursuant to this model, these capabilities are insufficient to deem the AI robot a perpetrator of an offense. These capabilities resemble the parallel capabilities of a mentally limited person, such as a child, or of a person who is mentally incompetent or who lacks a criminal state of mind.⁴¹ Legally, when an offense is committed by an innocent agent, like when a person causes a child,⁴² a person who is mentally incompetent⁴³ or who lacks a

³⁸ HALL, *supra* note 26.

³⁹ *Id.* at 105-45, 183.

⁴⁰ *See* discussion *supra* Part II.

⁴¹ HALL, *supra* note 26, at 232.

⁴² *See, e.g.,* Maxey v. United States, 30 App. D.C. 63 (1907); Commonwealth v. Hill, 11 Mass. 136 (1814); R. v. Michael, 169 Eng. Rep. 48 (1840).

criminal state of mind to commit an offense,⁴⁴ that person is criminally liable as a perpetrator-by-another.⁴⁵ In such cases, the intermediary is regarded as a mere instrument, albeit a sophisticated instrument, while the party orchestrating the offense (the perpetrator-by-another) is the real perpetrator as a principal in the first degree and is held accountable for the conduct of the innocent agent.⁴⁶ The perpetrator's liability is determined on the basis of that conduct⁴⁷ and his own mental state.⁴⁸

The derivative question relative to AI Robots is: Who is the perpetrator-by-another?

There are two candidates: the first is the programmer of the AI software installed in the specific robot and the second is the user. A programmer of AI software might design a program in order to commit offenses via the AI robot. For example, a programmer designs software for an operating robot. The robot is intended to be placed in a factory, and its software is designed to torch the factory at night when no one is there. The robot committed the arson, but the programmer is deemed the perpetrator. The second person who might be considered the perpetrator-by-another is the user of the AI robot. The user did not program the software, but he

⁴³ *See, e.g.*, Johnson v. State, 38 So. 182 (Ala. 1904); People v. Monks, 24 P.2d 508 (Cal. Dist. Ct. App. 1933).

⁴⁴ *See, e.g.*, United States v. Bryan, 483 F.2d 88 (3rd Cir. 1973); Boushea v. United States, 173 F.2d 131 (8th Cir. 1949); People v. Mutchler, 140 N.E. 820 (Ill. 1923); State v. Runkles, 605 A.2d 111 (Md. 1992); Parnell v. State, 912 S.W.2d 422 (Ark. 1996); State v. Thomas, 619 S.W.2d 513 (Tenn. 1981).

⁴⁵ *See generally* Morrisey v. State, 620 A.2d 207 (Del. 1993); State v. Fuller, 552 S.E.2d 282 (S.C. 2001); Gallimore v. Commonwealth, 436 S.E.2d 421 (Va. 1993).

⁴⁶ *Id.*

⁴⁷ *See generally* Dusenbery v. Commonwealth, 263 S.E.2d 392, 392 (Va. 1980).

⁴⁸ *See generally* United States v. Tobon-Builes, 706 F.2d 1092 (11th Cir. 1983); United States v. Ruffin, 613 F.2d 408 (2nd Cir. 1979).

uses the AI robot, including its software, for his own benefit. For example, a user purchases a servant-robot, which is designed to execute any order given by its master. The specific user is identified by the robot as that master, and the master orders the robot to assault any invader of the house. The robot executes the order exactly as ordered. This is not different than a person who orders his dog to attack any trespasser. The robot committed the assault, but the user is deemed the perpetrator.

In both scenarios, the actual offense was committed by the AI robot. The programmer or the user did not perform any action conforming to the definition of a specific offense; therefore, they do not meet the *actus reus* requirement of the specific offense. The perpetration-by-another liability model considers the action committed by the AI robot as if it had been the programmer's or the user's action. The legal basis for that is the instrumental usage of the AI robot as an innocent agent. No mental attribute required for the imposition of criminal liability is attributed to the AI robot.⁴⁹ When programmers or users use an AI robot instrumentally, the commission of an offense by the AI robot is attributed to them. The mental element required in the specific offense already exists in their minds.⁵⁰ The programmer had criminal intent when he ordered the commission of the arson, and the user had criminal intent when he ordered the commission of the assault, even though these offenses were actually committed through an AI robot.

This liability model does not attribute any mental capability, or any human mental capability, to the AI robot. According to this model, there is no legal difference between an AI robot and a screwdriver or an animal. When a burglar uses a screwdriver in order to open up a

⁴⁹ The AI robot is used as an instrument and not as a participant, although it uses its features of processing information. See George R. Cross & Cary G. Debessonnet, *An Artificial Intelligence Application in the Law: CCLIPS, A Computer Program that Processes Legal Information*, 1 HIGH TECH. L.J. 329, 362 (1986).

⁵⁰ HALL, *supra* note 26.

window, he uses the screwdriver instrumentally, and the screwdriver is not criminally liable. The screwdriver's "action" is, in fact, the burglar's. This is the same legal situation when using an animal instrumentally. An assault committed by a dog by order of its master is, in fact, an assault committed by the master.

This kind of legal model might be suitable for two types of scenarios. The first scenario is using an AI robot to commit an offense without using its advanced capabilities. The second scenario is using a very old version of an AI robot, which lacks the modern advanced capabilities of the modern AI robots. In both scenarios, the use of the AI robot is instrumental; it is usage of an AI robot, given its ability to execute an order, to commit an offense. A screwdriver cannot execute such an order; a dog can. A dog cannot execute complicated orders; an AI robot can.⁵¹

The perpetration-by-another liability model is not suitable when an AI robot decides to commit an offense based on its own accumulated experience or knowledge. This model is not suitable when the software of the AI robot was not designed to commit the specific offense, but was committed by the AI robot nonetheless. This model is also not suitable when the specific AI robot functions not as an innocent agent, but as a semi-innocent agent.⁵² However, the perpetration-by-another liability model might be suitable when a programmer or user makes instrumental usage of an AI robot, but without using the AI robot's advanced capabilities. The legal result of applying this model is that the programmer and the user are both criminally liable for the specific offense committed, while the AI robot has no criminal liability whatsoever.

⁵¹ Compare Andrew J. Wu, *From Video Games to Artificial Intelligence: Assigning Copyright Ownership to Works Generated by Increasingly Sophisticated Computer Programs*, 25 AIPLA Q.J. 131 (1997), with Timothy L. Butler, *Can a Computer be an Author: Copyright Aspects of Artificial Intelligence*, 4 HASTINGS COMM. & ENT. L.J. 707 (1982).

⁵² See generally NICOLA LACEY & CELIA WELLS, *RECONSTRUCTING CRIMINAL LAW: CRITICAL PERSPECTIVES ON CRIME AND THE CRIMINAL PROCESS* 53 (2d ed. 1998).

B. The Natural-Probable-Consequence Liability Model: Foreseeable Offenses Committed by AI Robots

The second model of criminal liability assumes deep involvement of the programmers or users in the AI robot's daily activities, but without any intention of committing any offense via the AI robot. For instance, one scenario would be when an AI robot commits an offense during the execution of its daily tasks. The programmers or users had no knowledge of the offense until it had already been committed. They did not plan to commit any offense, and they did not participate in any part of the commission of that specific offense.

An example of such a scenario is when an AI robot or software is designed to function as an automatic pilot. As part of the mission of flying the plane, the AI robot is programmed to protect the mission itself. During the flight, the human pilot activates the automatic pilot (which is the AI robot), and the program is initiated. At some point after activation of the automatic pilot, the human pilot sees an approaching storm and tries to abort the mission and return to base. The AI robot deems the human pilot's action as a threat to the mission and takes action in order to eliminate that threat; it may attempt to cut off the air supply to the pilot or activate the ejection seat. Whatever defense tactic is taken, the human pilot is killed as a result of the AI robot's actions. Obviously, the programmer had not intended to kill anyone, especially not the human pilot, but nonetheless, the human pilot was killed by the AI robot's programmed actions.

In this example, the first model is not legally suitable. The first model assumes *mens rea*, the criminal intent of the programmers or users to commit an offense via the instrumental use of some of the AI robot's capabilities.⁵³ This is not the legal situation in the case of the automatic pilot. In this case, the programmers or users had no knowledge of the committed offense; they had not planned it, and had not intended to commit the offense using the AI robot. For such

⁵³ See discussion *supra* Part III.A.

circumstances, the natural-probable-consequence liability model may create a more suitable legal response. This model is based upon the ability of the programmers or users to foresee the potential commission of offenses.

According to the second model, a person might be held accountable for an offense, if that offense is a natural and probable consequence of that person's conduct. Originally, the natural-probable-consequence liability was used to impose criminal liability upon accomplices, when one committed an offense, which had not been planned by all of them and which was not part of a conspiracy.⁵⁴ The established rule prescribed by courts and commentators is that accomplice liability extends to acts of a perpetrator that were a "natural and probable consequence"⁵⁵ of a criminal scheme that the accomplice encouraged or aided.⁵⁶ The natural-probable-consequence liability has been widely accepted in accomplice liability statutes and recodifications.⁵⁷

Natural-probable-consequence liability seems to be legally suitable for situations where an AI robot committed an offense, but the programmer or user had no knowledge of it, had not intended it and had not participated in it. The natural-probable-consequence liability model only requires the programmer or user to be in a mental state of negligence, not more. Programmers or users are not required to know about any forthcoming commission of an offense as a result of

⁵⁴ See generally *United States v. Powell*, 929 F.2d 724 (D.C. Cir. 1991).

⁵⁵ *Id.*

⁵⁶ See generally WILLIAM M. CLARK & WILLIAM L. MARSHALL, *LAW OF CRIMES* 529 (7th ed. 1967); Francis Bowes Sayre, *Criminal Responsibility for the Acts of Another*, 43 HARV. L. REV. 689 (1930); and see, e.g., *People v. Prettyman*, 926 P.2d 1013 (Cal. 1996); *Chance v. State*, 685 A.2d 351 (Del. 1996); *Ingram v. United States*, 592 A.2d 992 (D.C. 1991); *Richardson v. State*, 697 N.E.2d 462 (Ind. 1998); *Mitchell v. State*, 971 P.2d 813 (Nev. 1998); *State v. Carrasco*, 928 P.2d 939 (N.M. 1996); *State v. Jackson*, 976 P.2d 1229 (Wash. 1999).

⁵⁷ See, e.g., *United States v. Andrews*, 75 F.3d 552, 553-57 (9th Cir. 1996); *State v. Kaiser*, 918 P.2d 629, 632-39 (Kan. 1996).

their activity, but are required to know that such an offense is a natural, probable consequence of their actions.

A negligent person, in a criminal context, is a person who has no knowledge of the offense; rather, a reasonable person should have known about it since the specific offense is a natural probable consequence of that person's conduct.⁵⁸ Thus, the programmers or users of an AI robot, who should have known about the probability of the forthcoming commission of the specific offense, are criminally liable for the specific offense, even though they did not actually know about it. This is the fundamental legal basis for criminal liability in negligence cases.⁵⁹ Negligence is, in fact, an omission of awareness or knowledge.⁶⁰ The negligent person omitted knowledge, not acts.⁶¹

The natural-probable-consequence liability model would permit liability to be predicated upon negligence, even when the specific offense requires a different state of mind.⁶² This is not valid in relation to the person who personally committed the offense, but rather, is considered valid in relation to the person who was not the actual perpetrator of the offense, but was one of its intellectual perpetrators.⁶³ Reasonable programmers or users should have foreseen the

⁵⁸ See generally Robert P. Fine & Gary M. Cohen, *Is Criminal Negligence a Defensible Basis for Criminal Liability?*, 16 BUFF. L. REV. 749, 749-52 (1966); Herbert L.A. Hart, *Negligence, Mens rea and Criminal Responsibility*, in OXFORD ESSAYS IN JURISPRUDENCE 29 (1961); Donald Stuart, *Mens rea*, in NEGLIGENCE AND ATTEMPTS, 1968 CRIM. L. REV. 647 (1968).

⁵⁹ HALL, *supra* note 26, at 114-40.

⁶⁰ *Id.*

⁶¹ *Id.*

⁶² THE AMERICAN LAW INSTITUTE, MODEL PENAL CODE: OFFICIAL DRAFT AND EXPLANATORY NOTES § 2.05 (1962, 1985) [hereinafter Model Penal Code]; *and see, e.g.*, State v. Linscott, 520 A.2d 1067, 1069 (Me. 1987); People v. Luparello, 231 Cal. Rptr. 832 (Cal. Ct. App. 1987).

offense, and prevented it from being committed by the AI robot. However, the legal results of applying the natural-probable-consequence liability model to the programmer or user differ in two different types of factual cases. The first type of case is when the programmers or users were negligent while programming or using the AI robot but had no criminal intent to commit any offense. The second type of case is when the programmers or users programmed or used the AI robot knowingly and willfully in order to commit one offense via the AI robot, but the AI robot deviated from the plan and committed some other offense, in addition to or instead of the planned offense.

The first type of case is a pure case of negligence. The programmers or users acted or omitted negligently; therefore, there is no reason why they should not be held accountable for an offense of negligence, if there is such an offense in the specific legal system. Thus, as in the example above, where a programmer of an automatic pilot negligently programmed it to defend its mission with no restrictions on the taking of human life, the programmer is negligent and liable for the homicide of the human pilot. Consequently, if there is a specific offense of negligent homicide in that legal system, this is the most severe offense, for which the programmer might be held accountable because manslaughter or murder requires knowledge or intent.

The second type of case resembles the basic idea of the natural-probable-consequence liability in accomplice liability cases. The dangerousness of the very association or conspiracy whose aim is to commit an offense is the legal reason for more severe accountability to be imposed upon the cohorts. For example, a programmer programs an AI robot to commit a violent robbery of a bank, but the programmer did not program the AI robot to kill anyone.

⁶³ See sources cited *supra* note 62.

During the execution of the violent robbery, the AI robot kills one of the people present at the bank who resisted the robbery. In such cases, the criminal negligence liability alone is insufficient. The danger posed by such a situation far exceeds negligence.

As a result, according to the natural-probable-consequence liability model, when the programmers or users programmed or used the AI robot knowingly and willfully in order to commit one offense via the AI robot, but the AI robot deviated from the plan and committed another offense, in addition to or instead of the planned offense, the programmers or users shall be held accountable for the offense itself as if it had been committed knowingly and willfully. In the above example of the robbery, the programmer shall be held criminally accountable for the robbery (if committed), as well as for the killing as an offense of manslaughter or murder, which requires knowledge and intent.⁶⁴

The question still remains: What is the criminal liability of the AI robot itself when the natural-probable-consequence liability model is applied? In fact, there are two possible outcomes. If the AI robot acted as an innocent agent, without knowing anything about the criminal prohibition, it is not held criminally accountable for the offense it committed. Under such circumstances, the actions of the AI robot were not different from the actions of the AI robot under the first model (the perpetration-by-another liability model⁶⁵). However, if the AI robot did not act merely as an innocent agent, then the AI robot itself shall be held criminally liable for the specific offense directly, in addition to the criminal liability of the programmer or

⁶⁴ See, e.g., *United States v. Greer*, 467 F.2d 1064 (7th Cir. 1972); *People v. Cooper*, 743 N.E.2d 32 (Ill. 2000); *People v. Michalow*, 128 N.E. 228 (N.Y. 1920); *People v. Little*, 107 P.2d 634 (Cal. Dist. Ct. App. 1941); *People v. Cabalero*, 87 P.2d 364 (Cal. Dist. Ct. App. 1939); *People v. Weiss*, 9 N.Y.S.2d 1 (N.Y. App. 1939); *R v. Cunningham*, 3 W.L.R. 76 (1957); *R v. Faulkner*, 13 Cox C.C. 550 (1876).

⁶⁵ See *supra* Part III.A.

user pursuant to the natural-probable-consequence liability model. The direct liability model of AI robots is the third model, as described hereunder.

C. The Direct Liability Model: AI Robots as Direct Subjects of Criminal Liability

The third model does not assume any dependence of the AI robot on a specific programmer or user. The third model focuses on the AI robot itself.⁶⁶ As discussed above, criminal liability for a specific offense is mainly comprised of the factual element (*actus reus*) and the mental element (*mens rea*) of that offense.⁶⁷ Any person attributed with both elements of the specific offense is held criminally accountable for that specific offense.⁶⁸ No other criteria are required in order to impose criminal liability.⁶⁹ A person might possess further capabilities, but, in order to impose criminal liability, the existence of the factual element and the mental element required to impose liability for the specific offense is quite enough.⁷⁰ In order to impose criminal liability on any kind of entity, the existence of these elements in the specific entity must be proven.⁷¹ When it has been proven that a person committed the offense in question with

⁶⁶ See generally Steven J. Frank, *Tort Adjudication and the Emergence of Artificial Intelligence Software*, 21 SUFFOLK U. L. REV. 623 (1987); Sam N. Lehman-Wilzig, *Frankenstein Unbound: Towards a Legal Definition of Artificial Intelligence*, 13 FUTURES 442 (1981); Maruerite E. Gerstner, *Liability Issues with Artificial Intelligence Software*, 33 SANTA CLARA L. REV. 239 (1993); Richard E. Susskind, *Expert Systems in Law: A Jurisprudential Approach to Artificial Intelligence and Legal Reasoning*, 49 MOD. L. REV. 168 (1986).

⁶⁷ HALL, *supra* note 26.

⁶⁸ *Id.*

⁶⁹ *Id.*

⁷⁰ *Id.* at 185-86.

⁷¹ *Id.* at 183.

knowledge or intent, that person is held criminally liable for that offense.⁷² The relevant questions regarding the criminal liability of AI robots are: How can these robots fulfill the requirements of criminal liability? Do AI robots differ from humans in this context?

An AI algorithm might have many features and qualifications far exceeding those of an average human, but such features or qualifications are not required in order to impose criminal liability.⁷³ When a human or corporation fulfills the requirements of both the factual element and the mental element, criminal liability is imposed.⁷⁴ If an AI robot is capable of fulfilling the requirements of both the factual element and the mental element, and, in fact, it actually fulfills them, there is presumptively nothing to prevent criminal liability from being imposed on that AI robot.

Generally, the fulfillment of the factual element requirement of an offense is easily attributed to AI robots. As long as an AI robot controls a mechanical or other mechanism to move its moving parts, any act might be considered as performed by the AI robot. Thus, when an AI robot activates its electric or hydraulic arm and moves it, this might constitute an act. For example, in the specific offense of assault, such an electric or hydraulic movement of an AI robot that hits a person standing nearby is considered as fulfilling the *actus reus* requirement of the offense of assault. When an offense might be committed due to an omission, it is even simpler. Under this scenario, the AI robot is not required to act at all; its very inaction is the legal basis for criminal liability, as long as there was a duty to act. If a duty to act is imposed upon the AI

⁷² HALL, *supra* note 26, at 105-45, 183.

⁷³ *Id.* at 70-71; *see supra* Part II.

⁷⁴ *Id.*

robot and it fails to act, the *actus reus* requirement of the specific offense is fulfilled by way of an omission.

In most cases, the attribution of the mental element of offenses to AI robots is the real legal challenge. The attribution of the mental element differs from one AI technology to other. Most cognitive capabilities developed in modern AI technology are immaterial to the question of the imposition of criminal liability.⁷⁵ Creativity is a human feature that some animals have, but creativity is a not a requirement for imposing criminal liability.⁷⁶ Even the most uncreative persons are held criminally liable. The sole mental-state requirements to impose criminal liability are knowledge, intent, negligence, or the *mens rea* required in the specific offense and under the general theory of criminal law.⁷⁷

Knowledge is defined as sensory reception of factual data and the understanding of that data.⁷⁸ Most AI systems are well equipped for such reception because they possess sensory receptors for sights, voices, physical contact, touch, and the like.⁷⁹ These receptors transfer the

⁷⁵ HALL, *supra* note 26.

⁷⁶ *Id.*

⁷⁷ *Id.*

⁷⁸ *See generally* WILLIAM JAMES, *THE PRINCIPLES OF PSYCHOLOGY* (1890); HERMANN VON HELMHOLTZ, *THE FACTS OF PERCEPTION* (1878). In this context, knowledge and awareness are identical. *See, e.g.*, *United States v. Youts*, 229 F.3d 1312 (10th Cir. 2000); *United States v. Wert-Ruiz*, 228 F.3d 250 (3th Cir. 2000); *United States v. Ladish Malting Co.*, 135 F.3d 484 (7th Cir. 1998); *United States v. Spinney*, 65 F.3d 231 (1st Cir. 1995); *United States v. Jewell*, 532 F.2d 697 (9th Cir. 1976); *State v. Sargent*, 594 A.2d 401 (Vt. 1991); *State v. Wyatt*, 482 S.E.2d 147 (W. Va. 1996). *See also* Model Penal Code, *supra* note 62, at § 2.02(2)(b), which provides that, “A person acts *knowingly* with a respect to a material element of an offense when: (i) [if] he is *aware* that his conduct is of that nature or that such circumstances exist; and (ii) [if] he is *aware* that it is practically certain that his conduct will cause such a result” (emphasis added).

⁷⁹ *See generally* Margaret A. Boden, *Has AI Helped Psychology?*, in *THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE* 108 (Derek Partridge & Yorick Wilks eds., 2006); David Marr, *AI: A*

factual data received to central processing units that analyze the data.⁸⁰ The process of analysis in AI systems parallels that of human understanding.⁸¹ The human brain understands and analyzes the data received by eyes, ears, and hands.⁸² Advanced AI algorithms are trying to imitate human cognitive processes because these processes are not so different.⁸³

Specific intent is the strongest of the mental element requirements.⁸⁴ Specific intent is the existence of a purpose or an aim that a factual event will occur.⁸⁵ The specific intent required to establish liability for murder is a purpose or an aim that a certain person will die.⁸⁶ As a result of the existence of such intent, the perpetrator of the offense commits the offense, *i.e.*, he performs the factual element of the specific offense.⁸⁷ This situation is not unique to humans,

Personal View, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 97 (Derek Partridge & Yorick Wilks eds., 2006).

⁸⁰ See generally Derek Partridge, *What's in an AI Program?*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 112 (Derek Partridge & Yorick Wilks eds., 2006).

⁸¹ *Id.*

⁸² *Id.*

⁸³ See generally Daniel C. Dennett, *Evolution, Error, and Intentionality*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 190 (Derek Partridge & Yorick Wilks eds., 2006); B. Chandraswkar, *What Kind of Information Processing is Intelligence?*, in THE FOUNDATIONS OF ARTIFICIAL INTELLIGENCE 14 (Derek Partridge & Yorick Wilks eds., 2006).

⁸⁴ See generally Robert Batey, *Judicial Exploration of Mens rea Confusion at Common Law and Under the Model Penal Code*, 18 GA. ST. U. L. REV. 341 (2001).

⁸⁵ See, e.g., *Carter v. United States*, 530 U.S. 255 (2000); *United States v. Randolph*, 93 F.3d 656 (9th Cir. 1996); *United States v. Torres*, 977 F.2d 321 (7th Cir. 1992); *Frey v. United States*, 708 So.2d 918 (Fla. 1998); *State v. Neuzil*, 589 N.W.2d 708 (Iowa 1999); *State v. Daniels*, 109 So.2d 896 (La. 1958); *People v. Disimone*, 650 N.W.2d 436 (Mich. Ct. App. 2002); *People v. Henry*, 607 N.W.2d 767 (Mich. Ct. App. 1999).

⁸⁶ For the intent-to-kill murder, see WAYNE R. LAFAVE, CRIMINAL LAW 733-34 (4th ed. 2003).

⁸⁷ *Id.*

and an AI robot might be programmed to have a similar purpose or an aim and to take actions to achieve that purpose. This is specific intent.⁸⁸

One might assert that humans have feelings that cannot be imitated by AI robots, not even by the most advanced robots. Examples of such feelings are love, affection, hatred, or jealousy.⁸⁹ This theory might be correct in relation to the technology of the beginning of the 21st Century;⁹⁰ however, these feelings are rarely required in specific offenses. Most specific offenses are satisfied by knowledge of the existence of the factual element.⁹¹ Few offenses require specific intent in addition to knowledge.⁹² Almost all other offenses are satisfied by much less than that (*e.g.*, negligence, recklessness, strict liability). Perhaps in a few specific offenses that do require certain feelings (*e.g.*, crimes of racism, hate⁹³), criminal liability cannot be imposed upon AI robots, which have no such feelings, but in any other specific offense, it is not a barrier.⁹⁴

⁸⁸ LAFAVE, *supra* note 86.

⁸⁹ Capps, *supra* note 15.

⁹⁰ *Id.*

⁹¹ HALL, *supra* note 26, at 632.

⁹² *Id.*

⁹³ See generally Elizabeth A. Boyd, Richard A. Berk & Karl M. Hammer, "Motivated by Hatred or Prejudice": Categorization of Hate-Motivated Crimes in Two Police Divisions, 30 LAW & SOC'Y REV. 819, 842 (1996); Theresa Suozzi, F. Matt Jackson, Jeff Kauffman et al., *Crimes Motivated by Hatred: The Constitutionality and Impact of Hate Crimes Legislation in the United States*, 1 SYRACUSE J. LEGIS. & POL'Y 29 (1995).

⁹⁴ HALL, *supra* note 26, at 105-45.

If a person fulfills the requirements of both the factual element and the mental element of a specific offense, then the person is held criminally liable.⁹⁵ Why should an AI robot that fulfills all elements of an offense be exempt from criminal liability? One might argue that some segments of human society are exempt from criminal liability even if both the factual and mental elements have been established. Such segments of society are infants and the mentally ill.⁹⁶ A specific order in criminal law exempts infants from criminal liability.⁹⁷ The social rationale behind the infancy defense is to protect infants from the harmful consequences of the criminal process and to handle them in other social frameworks.⁹⁸ Do such frameworks exist for AI robots? The original legal rationale behind the infancy defense was the fact that infants are as

⁹⁵ HALL, *supra* note 26, at 70-211.

⁹⁶ *See supra* notes 40-45 and accompanying text.

⁹⁷ *Id.*; and *see, e.g.*, MINN. STAT. § 9913 (1927); MONT. REV. CODE § 10729 (1935); N.Y. PENAL CODE § 816 (1935); OKLA. STAT. § 152 (1937); UTAH REV. STAT. 103-i-40 (1933); *State v. George*, 54 A. 745 (Del. 1902); *Heilman v. Commonwealth*, 1 S.W. 731 (Ky. 1886); *State v. Aaron*, 4 N.J.L. 269 (N.J. 1818); *McCormack v. State*, 15 So. 438 (Ala. 1894); *Little v. State*, 554 S.W.2d 312 (Ark. 1977); *Clay v. State*, 196 So. 462 (Fla. 1940); *In re Devon T.*, 584 A.2d 1287 (Md. 1991); *State v. Dillon*, 471 P.2d 553 (Idaho 1970); *State v. Jackson*, 142 S.W.2d 45 (Mo. 1940).

⁹⁸ *See generally* Frederick J. Ludwig, *Rationale of Responsibility for Young Offenders*, 29 NEB. L. REV. 521 (1950); *In re Tyvonne*, 558 A.2d 661 (Conn. 1989); Andrew Walkover, *The Infancy Defense in the New Juvenile Court*, 31 UCLA L. REV. 503 (1984); Keith Foren, *In Re Tyvonne M. Revisited: The Criminal Infancy Defense in Connecticut*, 18 Q. L. REV. 733 (1999); Michael Tonry, *Rethinking Unthinkable Punishment Policies in America*, 46 UCLA L. REV. 1751 (1999); Andrew Ashworth, *Sentencing Young Offenders*, in PRINCIPLED SENTENCING: READINGS ON THEORY AND POLICY 294 (Andrew von Hirsch, Andrew Ashworth & Julian Roberts eds., 3d ed. 2009); Franklin E. Zimring, *Rationales for Distinctive Penal Policies for Youth Offenders*, in PRINCIPLED SENTENCING: READINGS ON THEORY AND POLICY 316 (Andrew von Hirsch, Andrew Ashworth & Julian Roberts eds., 3d ed. 2009); Andrew von Hirsch, *Reduced Penalties for Juveniles: The Normative Dimension*, in PRINCIPLED SENTENCING: READINGS ON THEORY AND POLICY 323 (Andrew von Hirsch, Andrew Ashworth & Julian Roberts eds., 3d ed. 2009).

yet incapable of comprehending what was wrong in their conduct (*doli incapax*).⁹⁹ Later, children can be held criminally liable if the presumption of mental incapacity was refuted by proof that the child was able to distinguish between right and wrong.¹⁰⁰ Could that be similarly applied to AI robots? Most AI algorithms are capable of analyzing permitted and forbidden.

The mentally ill are presumed to lack the fault element of the specific offense, due to their mental illness (*doli incapax*).¹⁰¹ The mentally ill are unable to distinguish between right and wrong (*cognitive capabilities*)¹⁰² and to control impulsive behavior.¹⁰³ When an AI algorithm functions properly, there is no reason for it not to use all of its capabilities to analyze the factual data received through its receptors. However, an interesting legal question would be whether a defense of insanity might be raised in relation to a malfunctioning AI algorithm, when its analytical capabilities become corrupted as a result of that malfunction.

⁹⁹ SIR EDWARD COKE, INSTITUTIONS OF THE LAWS OF ENGLAND: THIRD PART 4 (6th ed. 2001) (1681).

¹⁰⁰ MATTHEW HALE, HISTORIA PLACITORUM CORONAE 23, 26 (1736) [MATTHEW HALE, HISTORY OF THE PLEAS OF THE CROWN (1736)]; *and see, e.g.*, McCormack v. State, 15 So. 438 (Ala. 1894); Little v. State, 554 S.W.2d 312 (Ark. 1977); *In re Devon T.*, 584 A.2d 1287 (Md. 1991).

¹⁰¹ Benjamin B. Sendor, *Crime as Communication: An Interpretive Theory of the Insanity Defense and the Mental Elements of Crime*, 74 GEO. L.J. 1371, 1380 (1986); Joseph H. Rodriguez, Laura M. LeWinn & Michael L. Perlin, *The Insanity Defense Under Siege: Legislative Assaults and Legal Rejoinders*, 14 RUTGERS L.J. 397, 406-07 (1983); Homer D. Crotty, *The History of Insanity as a Defence to Crime in English Common Law*, 12 CAL. L. REV. 105 (1924).

¹⁰² *See generally* Edward de Grazia, *The Distinction of Being Mad*, 22 U. CHI. L. REV. 339 (1955); Warren P. Hill, *The Psychological Realism of Thurman Arnold*, 22 U. CHI. L. REV. 377 (1955); Manfred S. Guttmacher, *The Psychiatrist as an Expert Witness*, 22 U. CHI. L. REV. 325 (1955); Wilber G. Katz, *Law, Psychiatry, and Free Will*, 22 U. CHI. L. REV. 397 (1955); Jerome Hall, *Psychiatry and Criminal Responsibility*, 65 YALE L.J. 761 (1956).

¹⁰³ *See generally* John Barker Waite, *Irresistible Impulse and Criminal Liability*, 23 MICH. L. REV. 443, 454 (1925); Edward D. Hoedemaker, *"Irresistible Impulse" as a Defense in Criminal Law*, 23 WASH. L. REV. 1, 7 (1948).

When an AI robot establishes all elements of a specific offense, both factual and mental, it may be presumed that there is no reason to prevent imposition of criminal liability upon it for that offense. The criminal liability of an AI robot does not replace the criminal liability of the programmers or the users, if criminal liability is imposed on the programmers and users by any other legal path. Criminal liability is not to be divided, but rather, added. The criminal liability of the AI robot is imposed in addition to the criminal liability of the human programmer or user.

However, the criminal liability of an AI robot is not dependent upon the criminal liability of the programmer or user of that AI robot. As a result, if the specific AI robot was programmed or used by another AI robot, the criminal liability of the programmed or used AI robot is not influenced by that fact. The programmed or used AI robot shall be held criminally accountable for the specific offense pursuant to the direct liability model, unless it was an innocent agent. In addition, the programmer or user of the AI robot shall be held criminally accountable for that very offense pursuant to one of the three liability models, according to its specific role in the offense. The chain of criminal liability might continue, if more parties are involved, whether human or AI robots.

There is no reason to eliminate the criminal liability of an AI robot or of a human, which is based on complicity between them. An AI robot and a human might cooperate as joint perpetrators, as accessories and abettors, or the like; thus, the relevant criminal liability might be imposed on them accordingly. Since the factual and mental capabilities of an AI robot are sufficient to impose criminal liability—that is, if these capabilities satisfy the legal requirements of joint perpetrators, or of accessories and abettors—then the relevant criminal liability as joint perpetrators, accessories and abettors, or the like should be imposed irrespective of whether the offender is an AI robot or a human.

Not only positive factual and mental elements may be attributed to AI robots; rather, all relevant negative fault elements should be attributable to AI robots. Most of these elements are expressed by the general defenses in criminal law, *e.g.*, self-defense, necessity, duress, or intoxication. For some of these defenses (justifications),¹⁰⁴ there is no material difference between humans and AI robots since they relate to a specific situation (*in rem*), regardless of the identity of the offender. For example, an AI robot serving under the local police force is given an order to arrest a person illegally. If the order is not manifestly illegal, the executer of the order is not criminally liable.¹⁰⁵ In that case, there is no difference whether the executer is human or an AI robot.

For other defenses (excuses and exempts),¹⁰⁶ some applications should be adjusted. For example, the intoxication defense is applied when the offender is under the physical influence of an intoxicating substance (*e.g.*, alcohol or drugs). The influence of alcohol on an AI robot is minor, at most, but the influence of an electronic virus that is infecting the operating system of the AI robot might be considered parallel to the influence of intoxicating substances on humans.

¹⁰⁴ See generally JOHN C. SMITH, JUSTIFICATION AND EXCUSE IN THE CRIMINAL LAW (1989); Anthony M. Dillof, *Unraveling Unknowing Justification*, 77 NOTRE DAME L. REV. 1547 (2002); Kent Greenawalt, *Distinguishing Justifications from Excuses*, 49 LAW & CONTEMP. PROBS. 89 (1986); Kent Greenawalt, *The Perplexing Borders of Justification and Excuse*, 84 COLUM. L. REV. 949 (1984); Thomas Morawetz, *Reconstructing the Criminal Defenses: The Significance of Justification*, 77 J. CRIM. L. & CRIMINOLOGY 277 (1986); Paul H. Robinson, *A Theory of Justification: Societal Harm as a Prerequisite for Criminal Liability*, 23 UCLA L. REV. 266 (1975); Paul H. Robinson, *Testing Competing Theories of Justification*, 76 N.C. L. REV. 1095 (1998).

¹⁰⁵ See generally Michael A. Musmanno, *Are Subordinate Officials Penally Responsible for Obeying Superior Orders which Direct Commission of Crime?*, 67 DICK. L. REV. 221 (1963).

¹⁰⁶ See generally Peter Arenella, *Convicting the Morally Blameless: Reassessing the Relationship Between Legal and Moral Accountability*, 39 UCLA L. REV. 1511 (1992); Sanford H. Kadish, *Excusing Crime*, 75 CAL. L. REV. 257 (1987); Andrew E. Lelling, *A Psychological Critique of Character-Based Theories of Criminal Excuse*, 49 SYRACUSE L. REV. 35 (1998).

Some other factors might be considered as being parallel to insanity or loss of control. It may be concluded that the criminal liability of an AI robot, according to the direct liability model, is not different from the relevant criminal liability of a human. In some cases, some adjustments are necessary, but substantively, it is the very same criminal liability based upon the same elements and examined under the same light.

D. Hybrids: Coordinating the Models

The possible liability models described above are not alternative models.¹⁰⁷ These models might be applied in combination to create a full image of criminal liability in the specific context of AI robot involvement. None of the possible models is mutually exclusive. Thus, applying the second model is possible as a single model for the specific offense, and it is possible as one part of a combination of two of the legal models or of all three of them. When the AI robot plays the role of an innocent agent in the perpetration of a specific offense, and the programmer is the only person who directed that perpetration, the application of the perpetration-by-another model (the first liability model¹⁰⁸) is the most appropriate legal model for that situation. In that same situation, when the programmer is itself an AI robot (when an AI robot programs another AI robot to commit a specific offense), the direct liability model (the third liability model¹⁰⁹) is most appropriate to be applied to the criminal liability of the programmer of the AI robot. The third liability model in that situation is applied in addition to the first liability model, and not in lieu thereof. Thus, in such situations, the AI robot programmer shall be

¹⁰⁷ See discussion *supra* Parts III.A-C.

¹⁰⁸ See discussion *supra* Parts III.A.

¹⁰⁹ See discussion *supra* Parts III.C.

criminally liable, pursuant to a combination of the perpetration-by-another liability model and the direct liability model.¹¹⁰

If the AI robot plays the role of the physical perpetrator of the specific offense, but that very offense was not planned to be perpetrated, then the application of the natural-probable-consequence liability¹¹¹ model might be appropriate. The programmer might be deemed negligent if no offense had been deliberately planned to be perpetrated. Alternatively, the programmer might be held fully accountable for that specific offense if another offense had indeed been deliberately planned, but the specific offense that was perpetrated had not been part of the original criminal scheme. Nevertheless, when the programmer is not human, the direct liability model must be applied in addition to the simultaneous application of the natural-probable-consequence liability model; likewise, when the physical perpetrator is human while the planner is an AI robot.¹¹²

Hybrids of all three liability models create an opaque net of criminal liability. The combined and coordinated application of these three models reveals a new legal situation in the specific context of AI robots and criminal law. As a result, when AI robots and humans are involved, directly or indirectly, in the perpetration of a specific offense, it will be far more difficult to evade criminal liability. The social benefit to be derived from such a legal policy is of substantial value. All entities—human, legal or AI—become subject to criminal law. If the clearest purpose of the imposition of criminal liability is the application of legal social control in

¹¹⁰ See *supra* Parts III.A and III.C.

¹¹¹ See discussion *supra* Parts III.B.

¹¹² See *supra* Parts III.B and III.C.

the specific society, then the coordinated application of all three models is necessary in the very context of AI robots.

IV. GENERAL PUNISHMENT ADJUSTMENT CONSIDERATIONS

Let us assume an AI robot is criminally liable. Let us assume it is indicted, tried and convicted. After the conviction, the court is supposed to sentence that AI robot. If the most appropriate punishment under the specific circumstances is one year of imprisonment, for example, how can an AI robot practically serve such a sentence? How can capital punishment, probation or even a fine be imposed on an AI robot? What is the practical meaning of imprisonment? Where no bank account is available for the sentenced AI robot, what is the practical significance of fining it?

Similar legal problems have been raised when the criminal liability of corporations was recognized.¹¹³ Some asked how any of the legitimate penalties imposed upon humans could be applicable to corporations.¹¹⁴ The answer was simple and legally applicable.¹¹⁵ When a punishment can be imposed on a corporation as it is on humans, it is imposed without change.¹¹⁶ When the court adjudicates a fine, the corporation pays the fine in the same way that a human

¹¹³ See generally Gerard E. Lynch, *The Role of Criminal Law in Policing Corporate Misconduct*, 60 LAW & CONTEMP. PROBS. 23 (1997); Richard Gruner, *To Let the Punishment Fit the Organization: Sanctioning Corporate Offenders Through Corporate Probation*, 16 AM. J. CRIM. L. 1 (1988); Steven Walt & William S. Laufer, *Why Personhood Doesn't Matter: Corporate Criminal Liability and Sanctions*, 18 AM. J. CRIM. L. 263 (1991); John C. Coffee, Jr., "No Soul to Damn: No Body to Kick": An Unscandalised Inquiry Into the Problem of Corporate Punishment, 79 MICH. L. REV. 386 (1981); STEVEN BOX, POWER, CRIME AND MYSTIFICATION 16-79 (1983); Brent Fisse & John Braithwaite, *The Allocation of Responsibility for Corporate Crime: Individualism, Collectivism and Accountability*, 11 SYDNEY L. REV. 468 (1988).

¹¹⁴ *Id.*

¹¹⁵ *Id.*

¹¹⁶ *Id.*

pays the fine and in the same way that a corporation pays its bills in a civil context.¹¹⁷ However, when punishment of a corporation cannot be carried out in the same way as with humans, an adjustment is required.¹¹⁸ Such is the legal situation vis-à-vis AI robots.

The punishment adjustment considerations examine the theoretical foundations of any applied punishment. These considerations are applied in a similar manner and are comprised of three stages. Each stage may be explained by a question: (1) What is the fundamental significance of the specific punishment for a human?; (2) How does that punishment affect AI robots?; and (3) What practical punishments may achieve the same significance when imposed on AI robots? The most significant advantage of these punishment adjustment considerations is that the significance of the specific punishment remains identical when imposed on humans and AI robots. This method of punishment adjustment considerations is referred to below in some of the punishments used in modern societies, *e.g.*, capital punishment, imprisonment, suspended sentencing, community service and fines.

Capital punishment is considered the most severe punishment for humans, and there is no consensus regarding its constitutionality among the various jurisdictions.¹¹⁹ Capital punishment is the most effective method of incapacitating offenders as it relates to recidivism since, once the death sentence is carried out, the offender is obviously incapable of committing any further

¹¹⁷ See sources cited *supra* note 113.

¹¹⁸ *Id.*

¹¹⁹ See, *e.g.*, GG art. 102 (for the abolition of capital penalty in Germany in 1949); Murder (Abolition of Death Penalty) Act, 1965, 13-14 Eliz. 2, c. 71 (for murder in Britain in 1965); and for the debate in the United States, *e.g.*, *Wilkerson v. Utah*, 99 U.S. 130 (1878); *In re Kemmler*, 136 U.S. 436 (1890); *Gregg v. Georgia*, 428 U.S. 153 (1979); *Hunt v. Nuth*, 57 F.3d 1327 (4th Cir. 1995); *Campbell v. Wood*, 18 F.3d 662 (9th Cir. 1994); *Gray v. Lucas*, 710 F.2d 1048 (5th Cir. 1983); *People v. Daugherty*, 256 P.2d 911 (Cal. 1953); *Provenzano v. Moore*, 744 So.2d 413 (Fla. 1999); *Dutton v. State*, 91 A. 417 (Md. 1914).

offense.¹²⁰ The significance of capital punishment for humans is the deprivation of life.¹²¹ The “life” of an AI robot is its independent existence as an entity. Considering capital punishment’s efficacy in incapacitating offenders, the practical action that may achieve the same results as capital punishment when imposed on an AI robot is deletion of the AI software controlling the AI robot. Once the deletion sentence is carried out, the offending AI robot is incapable of committing any further offenses. The deletion eradicates the independent existence of the AI robot and is tantamount to the death penalty.

Imprisonment is one of the most popular sentences imposed in western legal systems for serious crimes.¹²² The significance of imprisonment for humans is the deprivation of human liberty and the imposition of severe limitations on human free behavior, freedom of movement and freedom to manage one’s personal life.¹²³ The “liberty” or “freedom” of an AI robot

¹²⁰ See generally ROBERT M. BOHM, DEATHQUEST: AN INTRODUCTION TO THE THEORY AND PRACTICE OF CAPITAL PUNISHMENT IN THE UNITED STATES 74-78 (1999); Austin Sarat, *The Cultural Life of Capital Punishment: Responsibility and Representation in ‘Dead Man Walking’ and ‘Last Dance’*, in THE KILLING STATE: CAPITAL PUNISHMENT IN LAW, POLITICS, AND CULTURE 226 (Austin Sarat ed., 1999); Peter Fitzpatrick, “Always More to Do”: *Capital Punishment and the (De)Composition of Law*, in THE KILLING STATE: CAPITAL PUNISHMENT IN LAW, POLITICS, AND CULTURE 117 (Austin Sarat ed., 1999).

¹²¹ See generally Franklin E. Zimring, *The Executioner’s Dissonant Song: On Capital Punishment and American Legal Values*, in THE KILLING STATE: CAPITAL PUNISHMENT IN LAW, POLITICS, AND CULTURE 137 (Austin Sarat ed., 1999); Anthony G. Amsterdam, *Selling a Quick Fix for Boot Hill: The Myth of Justice Delayed in Death Cases*, in THE KILLING STATE: CAPITAL PUNISHMENT IN LAW, POLITICS, AND CULTURE 148 (Austin Sarat ed., 1999).

¹²² See generally David J. Rothman, *For the Good of All: The Progressive Tradition in Prison Reform*, in HISTORY AND CRIME 271 (James A. Inciardi & Charles E. Faupel eds., 1980); MICHAEL WELCH, IRONIES OF IMPRISONMENT (2004); Roy D. King, *The Rise and Rise of Supermax: An American Solution in Search of a Problem?*, 1 PUNISHMENT & SOC’Y 163 (1999); CHASE RIVELAND, SUPERMAX PRISONS: OVERVIEW AND GENERAL CONSIDERATIONS (1999); JAMIE FELLNER & JOANNE MARINER, COLD STORAGE: SUPER-MAXIMUM SECURITY CONFINEMENT IN INDIANA (1997).

includes the freedom to act as an AI robot in the relevant area. For example, an AI robot in medical service has the freedom to participate in surgeries, or an AI robot in a factory has the freedom to manufacture. Considering the nature of a sentence of imprisonment, the practical action that may achieve the same effects as imprisonment when imposed on an AI robot is to put the AI robot out of use for a determinate period. During that period, no action relating to the AI robot's freedom is allowed, and thus its freedom or liberty is restricted.

Suspended sentencing is a very popular intermediate sanction in western legal systems for increasing the deterrent effect on offenders in lieu of actual imprisonment.¹²⁴ The significance of a suspended sentence for humans is the very threat of imprisonment if the human commits a specific offense or a type of specific offense.¹²⁵ If the human commits such an offense, a sentence of imprisonment will be imposed for the first offense in addition to the sentencing for the second offense.¹²⁶ As a result, humans are deterred from committing another offense and from becoming a recidivist offender.¹²⁷ Practically, a suspended sentence is imposed

¹²³ See generally Richard Korn, *The Effects of Confinement in the High Security Unit in Lexington*, 15 SOC. JUST. 8 (1988); Holly A. Miller, *Reexamining Psychological Distress in the Current Conditions of Segregation*, 1 J. CORRECTIONAL HEALTH CARE 39 (1994); FRIEDA BERNSTEIN, *THE PERCEPTION OF CHARACTERISTICS OF TOTAL INSTITUTIONS AND THEIR EFFECT ON SOCIALIZATION* (1979); BRUNO BETTELHEIM, *THE INFORMED HEART: AUTONOMY IN A MASS AGE* (1960); Marek M. Kaminski, *Games Prisoners Play: Allocation of Social Roles in a Total Institution*, 15 RATIONALITY & SOC'Y 188 (2003); JOHN IRWIN, *PRISONS IN TURMOIL* (1980); ANTHONY J. MANOCCHIO AND JIMMY DUNN, *THE TIME GAME: TWO VIEWS OF A PRISON* (1982).

¹²⁴ See generally MARC ANCEL, *SUSPENDED SENTENCE* (1971); Marc Ancel, *The System of Conditional Sentence or Sursis*, 80 L. Q. REV. 334 (1964).

¹²⁵ *Id.*

¹²⁶ *Id.*

¹²⁷ Anthony E. Bottoms, *The Suspended Sentence in England 1967-1978*, 21 BRITISH J. CRIMINOLOGY 1, 2-3 (1981).

only in the legal records.¹²⁸ No physical action is taken when a suspended sentence is imposed.¹²⁹ As a result, there is no difference between humans and AI robots. The statutory criminal records of the state do not differentiate between a suspended sentence imposed on humans, and those imposed on corporations or AI robots, as long as the relevant entity may be identified specifically and accurately.

Community service is also a very popular intermediate sanction in western legal systems in lieu of actual imprisonment.¹³⁰ In most legal systems, community service is a substitute for short sentences of actual imprisonment.¹³¹ In some legal systems, community service is imposed coupled with probation so that the offender “pays a price” for the damages he caused by committing the specific offense.¹³² The significance of community service for humans is compulsory contribution of labor to the community.¹³³ As discussed above, an AI robot can be engaged as a worker in very many areas.¹³⁴ When an AI robot works in a factory, its work is done for the benefit of the factory owners or for the benefit of the other workers in order to ease

¹²⁸ Bottoms, *supra* note 127.

¹²⁹ *Id.*

¹³⁰ *See generally* John Harding, *The Development of the Community Service*, in ALTERNATIVE STRATEGIES FOR COPING WITH CRIME 164 (Norman Tutt ed., 1978); HOME OFFICE, REVIEW OF CRIMINAL JUSTICE POLICY (1977); Andrew Willis, *Community Service as an Alternative to Imprisonment: A Cautionary View*, 24 PROBATION J. 120 (1977).

¹³¹ *Id.*

¹³² *See generally* Julie Leibrich, Burt Galaway & Yvonne Underhill, *Community Sentencing in New Zealand: A Survey of Users*, 50 FED. PROBATION 55 (1986); James Austin & Barry Krisberg, *The Unmet Promise of Alternatives*, 28 J. RES. IN CRIME & DELINQ. 374 (1982); Mark S. Umbreit, *Community Service Sentencing: Jail Alternatives or Added Sanction?*, 45 FED. PROBATION 3 (1981).

¹³³ *Id.*

¹³⁴ *See supra* p. 32.

and facilitate their professional tasks. In the same way that an AI robot works for the benefit of private individuals, it may work for the benefit of the community. When work for the benefit of the community is imposed on an AI robot as a compulsory contribution of labor to the community, it may be considered community service. Thus, the significance of community service is identical, whether imposed on humans or AI robots.

The adjudication of a fine is the most popular intermediate sanction in western legal systems in lieu of actual imprisonment.¹³⁵ The significance of paying a fine for humans is deprivation of some of their property, whether the property is money (a fine) or other property (forfeiture).¹³⁶ When a person fails to pay a fine, or has insufficient property to pay the fine, substitute penalties are imposed on the offender, particularly imprisonment.¹³⁷ The imposition of a fine on a corporation is identical to the imposition of a fine on a person, since both people and corporations have property and bank accounts. Thus, the payment of a fine is identical whether the paying entity is human or a corporate entity. However, most AI robots have no money or

¹³⁵ See generally GERHARDT GREBING, *THE FINE IN COMPARATIVE LAW: A SURVEY OF 21 COUNTRIES* (1982); NIGEL WALKER AND NICOLA PADFIELD, *SENTENCING: THEORY, LAW AND PRACTICE* (1996); Manfred Zuleeg, *Criminal Sanctions to be Imposed on Individuals as Enforcement Instruments in European Competition Law*, in *EUROPEAN COMPETITION LAW ANNUAL 2001: EFFECTIVE PRIVATE ENFORCEMENT OF EC ANTITRUST LAW 451* (Claus-Dieter Ehlermann & Isabela Atanasiu eds., 2001); Judith A. Greene, *Structuring Criminal Fines: Making an 'Intermediate Penalty' More Useful and Equitable*, 13 *JUST. SYS. J.* 37 (1988); Manfred Zuleeg, *Criminal Sanctions to be Imposed on Individuals as Enforcement Instruments in European Competition Law*, in *EUROPEAN COMPETITION LAW ANNUAL 2001: EFFECTIVE PRIVATE ENFORCEMENT OF EC ANTITRUST LAW 451* (Claus-Dieter Ehlermann & Isabela Atanasiu eds., 2001).

¹³⁶ See generally DOUGLAS C. McDONALD, JUDITH A. GREENE & CHARLES WORZELLA, *DAY-FINES IN AMERICAN COURTS: THE STATEN-ISLAND AND MILWAUKEE EXPERIMENTS* (1992); STEVE UGLOW, *CRIMINAL JUSTICE* (1995).

¹³⁷ See generally *Use of Short Sentences of Imprisonment by the Court*, REPORT OF THE SCOTTISH ADVISORY COUNCIL ON THE TREATMENT OF OFFENDERS (1960); FIORI RINALDI, *IMPRISONMENT FOR NON-PAYMENT OF FINES* (1976).

property of their own, nor have they any bank accounts. In effect, the imposition of fines on AI robots may be problematic.

Assuming, however, if an AI robot did have its own property or money, the imposition of a fine on it would be identical to the imposition of a fine on humans or corporations. For most humans and corporations, property is gained through labor.¹³⁸ When paying a fine, such property resulting from labor is transferred to the state.¹³⁹ That labor might be transferred to the state in the form of property or directly as labor. As a result, a fine imposed on an AI robot might be collected as money or property and as labor for the benefit of the community. When the fine is collected in the form of labor for the benefit of the community, it is not different from community service as described above.¹⁴⁰

Most common punishments are applicable to AI robots. The imposition of specific penalties on AI robots does not negate the nature of these penalties in comparison with their imposition on humans. Of course, some general punishment adjustment considerations are necessary in order to apply these penalties, but still, the nature of these penalties remains the same relative to humans and to AI robots.

V. CONCLUSION

If all of its specific requirements are met, criminal liability may be imposed upon any entity—human, corporate or AI robot. Modern times warrant modern legal measures in order to resolve today's legal problems. The rapid development of Artificial Intelligence technology requires current legal solutions in order to protect society from possible dangers inherent in

¹³⁸ JOHN LOCKE, TWO TREATISES OF GOVERNMENT (1689).

¹³⁹ *See supra* note 135.

¹⁴⁰ *See supra* pp. 33-34.

technologies not subject to the law, especially criminal law. Criminal law has a very important social function—that of preserving social order for the benefit and welfare of society. The threats upon that social order may be posed by humans, corporations or AI robots.

Traditionally, humans have been subject to criminal law, except when otherwise decided by international consensus. Thus, minors and mentally ill persons are not subject to criminal law in most legal systems around the world.¹⁴¹ Although corporations in their modern form have existed since the 14th Century,¹⁴² it took hundreds of years to subordinate corporations to the law, especially to criminal law.¹⁴³ For hundreds of years, the law stated that corporations are not subject to criminal law, as inspired by Roman law (*societas delinquere non potest*).¹⁴⁴ It was only in 1635 that an English court dared to impose criminal liability on a corporation.¹⁴⁵ Corporations participate fully in human life, and it was outrageous not to subject them to human laws since offenses are committed by corporations or through them. But, corporations have neither body nor soul. Legal solutions were developed so that in relation to criminal liability, they would be deemed capable of fulfilling all requirements of criminal liability, including

¹⁴¹ See discussion *supra* Part III.A.

¹⁴² WILLIAM SEARLE HOLDSWORTH, A HISTORY OF ENGLISH LAW 471-76 (1923).

¹⁴³ *Id.*

¹⁴⁴ See generally William Searle Holdsworth, *English Corporation Law in the 16th and 17th Centuries*, 31 YALE L.J. 382 (1922); WILLIAM ROBERT SCOTT, THE CONSTITUTION AND FINANCE OF ENGLISH, SCOTISH AND IRISH JOINT-STOCK COMPANIES TO 1720 462 (1912); BISHOP CARLETON HUNT, THE DEVELOPMENT OF THE BUSINESS CORPORATION IN ENGLAND 1800-1867 6 (Harvard University Press 1963).

¹⁴⁵ See, e.g., Case of Langforth Bridge, 79 Eng. Rep. 919 (K.B. 1635); R v. Inhabitants of Clifton, 101 Eng. Rep. 280 (K.B. 1794); R v. Inhabitants of Great Broughton, 98 Eng. Rep. 418 (K.B. 1771); R v. Mayor of Stratford upon Avon, 104 Eng. Rep. 636 (K.B. 1811); R v. The Mayor of Liverpool, 102 Eng. Rep. 529 (K.B. 1802); R v. Saintiff, 87 Eng. Rep. 1002 (K.B. 1705).

factual and mental elements.¹⁴⁶ These solutions were embodied in models of criminal liability and general punishment adjustment considerations.¹⁴⁷ It worked. In fact, it is still working, and very successfully.¹⁴⁸

Why should AI robots be different from corporations? AI robots are taking larger and larger parts in human activities, as do corporations. Offenses have already been committed by AI robots or through them. AI robots have no soul. Thus, there is no substantive legal difference between the idea of criminal liability imposed on corporations and on AI robots. It would be outrageous not to subordinate them to human laws, as corporations have been. As proposed by this article, models of criminal liability and general paths to impose punishment do exist. What else is needed?

¹⁴⁶ See generally Frederick Pollock, *Has the Common Law Received the Fiction Theory of Corporations?*, 27 L. Q. REV. 219 (1911).

¹⁴⁷ *Id.*

¹⁴⁸ *Id.*