

6-7-2019

A mesterséges általános intelligencia veszélyei és ígéretei

Brian S. Haney

Kövesse ezt és további munkáit a következő címen: <https://scholarship.law.nd.edu/jleg>

 A [Computational Engineering Commons](#), a [Computer Engineering Commons](#), a [Legislation Commons](#) és a [Science and Technology Policy Commons](#) része.

Ajánlott idézet

Brian S. Haney, *The Perils and Promises of Artificial General Intelligence*, J45. Legis. 151 (2018).

Elérhető a következő címen: <https://scholarship.law.nd.edu/jleg/vol45/iss2/1>

Ezt a cikket a Journal of Legislation at NDLScholarship ingyenesen és nyíltan hozzáférhetővé teszi. A Journal of Legislation című folyóiratba való felvételre az NDLScholarship egyik felhatalmazott szerkesztője fogadta el. További információért kérjük, forduljon a lawdr@nd.edu címre.

A MESTERSÉGES ÁLTALÁNOS INTELLIGENCIA VESZÉLYEI ÉS ÍGÉRETEI

Brian S. Haney[†]

BEVEZETÉS

A legtöbb ember úgy gondolja, hogy az ember és a technológia összeolvadása a távoli jövőben történhet meg; az igazság az, hogy az emberek már most is kiborgok. Egy okostelefonnal az ember gyakorlatilag bármilyen kérdésre gyorsan válaszolhat, korlátlan mennyiségű információt tárolhat a memóriájában, és bármilyen számítást elvégezhet.¹ A modern technológiai vállalatok az okostelefonokról adatokat gyűjtenek az emberekről, és közvetlenül a fejlett mesterséges intelligencia ("AI") rendszereken keresztül táplálják azokat.² A mesterséges intelligencia-rendszerek eleve maximalizálják a fogyasztók limbikus rendszerébe, az agy jutalomközpontjába érkező elektromos impulzusokat a gazdasági növekedés és fejlődés ösztönzése érdekében.³ Elon Musk a Nemzeti Kormányzók Szövetségének 2017-es nyári ülésén kijelentette: "[a] legnagyobb kockázat, amellyel civilizációnknak szembe kell néznie, a mesterséges intelligencia".⁴

Musk nincs egyedül, sőt, egyre több tudós és iparági vezető hívja fel a figyelmet arra, hogy a mesterséges intelligencia milyen egzisztenciális veszélyeket jelent az emberre nézve.⁵ A mesterséges intelligencia témájával foglalkozó jogi kutatások azonban vagy tagadták, vagy viszonylag figyelmen kívül hagyták a mesterséges intelligencia gyorsuló ütemű fejlődését.⁶ Ehelyett a mesterséges intelligencia szabályozásának szentelt jelenlegi jogi tudományosság arra ösztönözte a szabályozókat, hogy ne hagyják magukat megzavarni az "AI apokalipszisre" vonatkozó állításoktól, és erőfeszítéseiket a "közvetlenebb károokra" összpontosítsák.⁷ Összefoglalva, az AI-szabályozással foglalkozó jogi ösztöndíj messze elmaradt, és téves tanácsokat ad a szabályozóknak és a tudósoknak.⁸ Valójában minden olyan feladat, amelynek elvégzésére az emberek intelligenciát használnak, az AI automatizálásának célpontja.⁹ Továbbá gyakran előfordul, hogy amint egy mesterséges intelligenciával működő rendszer eléri az emberi szintű teljesítményt egy adott feladatban, röviddel ezután ugyanez a mesterséges intelligenciával működő rendszer meghaladja a legképzettebbek teljesítményét.

[†] Brian S. Haney, J.D. Notre Dame Law School 2018, a Martian Technologies vezérigazgatója. Külön köszönet illeti LaDarien Harrist, Mike Gallagher, Ryan Claudeanost, Delaney Foremant, Maria Munoz-Roblest, Brian Wongchaowart, Bill Greent, Ned Rooneyt és a Jogászujságot.

¹ *The Joe Rogan Experience #1169 - Elon Musk*, THE JOE ROGAN EXPERIENCE (2018. szeptember 6.), <https://www.youtube.com/watch?v=ycPr5-27vSI>.

² *Id.*

³ *Id.*

⁴ *Elon Musk at the National Governors Association 2017 Summer Meeting*, C-SPAN, 50:00 (July 15, 2017), <https://www.c-span.org/video/?c4676772/elon-musk-national-governors-association-2017-summer-meeting>.

⁵ *Lásd*: MILES BRUNDAGE ET AL., A MŰVES INTELLIGENCIA MÉLTALÁNOS HASZNÁLATA: FORECASTING, PREVENTION, AND MITIGATION (122018); *lásd még* MAX TEGMARK, LIFE BEING3.0 HUMAN IN THE AGE OF ARTIFICIAL INTELLIGENCE (2017).

⁶ *Lásd* Ryan Calo, *Mesterséges intelligencia politika*: DAVIS L. REV. 399,432 (2017).

⁷ *Lásd id.* 431.

⁸ *Lásd id.*

⁹ BRUNDAGE ET AL., *Supra* note. 5.

embereket a feladat elvégzésében.¹⁰ Sok AI-kutató arra számít, hogy az AI-rendszerek idővel minden feladatban elérik, majd meghaladják az emberi szintű teljesítményt.¹¹

A mesterséges intelligencia technológiája olyan jövőt alakít ki, ahol a hamis képek és videók olcsók, széles körben elérhetők és megkülönböztethetetlenek a valótól, ami teljesen átformálja azt a módot, ahogyan az emberek az igazságot a bizonyítékokkal társítják.¹² Még azok is, akik kételkednek abban, hogy a jövőben létre fog jönni a mesterséges általános intelligencia ("AGI"), azaz a bármilyen cél elérésére képes mesterséges intelligencia¹³, egyetértenek abban, hogy a mesterséges intelligenciának mélyreható következményei lesznek minden területen, többek között az egészségügyben, a jogban és a nemzetbiztonságban.¹⁴ E cikk célja kettős. Először is, ez a cikk meghatározza és elmagyarázza a mesterséges intelligencia csúcstechnológiáját, különös tekintettel a mély megerősítő tanulásra, a Google által 2013-ban kifejlesztett, áttörést jelentő gépi tanulás egyik típusára.¹⁵ Másodszor, ez a cikk három akadályt határoz meg, amelyeket a szabályozóknak le kell küzdeniük a mesterséges intelligencia szabályozása során.

Ez a cikk három fő módon járul hozzá a jelenlegi jogi és mesterséges intelligencia tudományhoz. Ez az első, amely a mély megerősítő tanulásra összpontosít, különösen a mesterséges intelligencia által jelentett egzisztenciális fenyegetésekre, és ez az első, amely a mesterséges intelligencia alapjául szolgáló formális modellekként foglalkozik. A cikk három részből áll. Az I. rész elmagyarázza a mesterséges intelligencia alapvető kifejezéseit és fogalmait, és a mesterséges intelligencia számos gyakorlati alkalmazását vizsgálja a modern iparban. A II. rész a mély megerősítő tanulást ismerteti, amely a mesterséges intelligencia viszonylag új keletű áttörése, és amely sok tudós szerint az AGI felé vezető utat jelent. A III. rész a mesterséges intelligencia szabályozásának témájával foglalkozó jogi tudományosságot vizsgálja, és három olyan kérdést tárgyal, amelyekkel a szabályozóknak foglalkozniuk kell ahhoz, hogy a mesterséges intelligencia erős szabályozási keretét kialakítsák.

I. MESTERSÉGES INTELLIGENCIA

A kortárs tudósok a mesterséges intelligencia számos különböző definícióját mutatták be. Max Tegmark, az MIT professzora például tömören úgy határozza meg az AI-t, mint "nem biológiai intelligencia".¹⁶ Ray Kurzweil, a Google munkatársa a mesterséges intelligenciát "olyan gépek létrehozásának művészeteként írta le, amelyek olyan funkciókat látnak el, amelyek elvégzéséhez embereknek intelligenciára van szükségük".¹⁷ Nils Nilsson, a Stanford professzora szerint az AI "a műtárgyak intelligens viselkedésével foglalkozik".¹⁸ Általánosságban és e cikk alkalmazásában az AI olyan intelligens gépek tanulmányozására és fejlesztésére utal, amelyek képesek az emberi kognitív funkciók - például jóslatok készítése, beszéd folyamatok vagy játék - gondolkodási folyamatait leutánozni.

Bár a mesterséges intelligencia különböző kategóriákat foglal magában, a mesterséges intelligencia szabályozásának összefüggésében két típusa a legfontosabb. Az első a szűk mesterséges intelligencia, más néven gyenge mesterséges intelligencia.¹⁹ A szűk értelemben vett mesterséges intelligencia korlátozott célok²⁰ elérésére képes, és az emberi intelligencia fejlesztésére irányuló kísérletekkel függ össze, szemben az emberi intelligencia megkettőzésével.

10 *Id.* 16.

11 *Id.*

12 GREG ALLEN & TANIEL CHAN, ARTIFICIAL INTELLIGENCE AND NATIONAL SECURITY (312017).

13 TEGMARK, 5,31. o. (2017).

14 ALLEN & CHAN, *fenti* megjegyzés 12.

15 Methods and Apparatus for Reinforcement Learning, U.S. Patent Application No. 14/097,862 (benyújtott 2013. december 5-én) (elérhető a <https://patents.google.com/patent/US20150100530A1/en> oldalon); *lásd még* TEGMARK, *su- pra* note at 5,84.

16 TEGMARK, *Supra* note5 .

17 RAY KURZWEIL, AZ INTELLIGENS GÉPEK KORA (141992).

18 NILS J. NILSSON, MESTERSÉGES INTELLIGENCIA: (11998).

19 NILS J. NILSSON, A MESTERSÉGES INTELLIGENCIA KERESÉSE (3882010).

20 *Lásd* TEGMARK, *fenti* megjegyzés 5.

emberi intelligencia.²¹ A mesterséges intelligencia második típusa a mesterséges általános intelligencia ("AGI"), más néven erős mesterséges intelligencia.²² Az AGI bizonyításához egy mesterséges intelligencia-ügynöknek képesnek kell lennie bármilyen cél elérésére.²³ Az AGI-hez kapcsolódik az az állítás, hogy egy programozott számítógép lehet elme, és legalább olyan jól tud gondolkodni, mint az emberek.²⁴ Végső soron az AGI számos AI-kutató jelenlegi célja.²⁵ Például az OpenAI, egy nonprofit szervezet, amely a terület úttörő kutatásait finanszírozza, a honlapján azt állítja, hogy küldetése "[d]járnyékolja és megvalósítja a biztonságos mesterséges általános intelligenciához vezető utat".²⁶ Mégis, egyelőre úgy tűnik, hogy csak szűk körű mesterséges intelligenciát fejlesztettek ki és alkalmaznak sikeresen.²⁷

A. AI A MODERN SZAKMAI IPARÁGAKBAN

A szűk értelemben vett mesterséges intelligencia alkalmazása világszerte felforgatja a modern iparágakat.²⁸ Még a jogi iparág sem mentesül ettől a maró erőttől.²⁹ A technológiával támogatott felülvizsgálat ("TAR") forradalmasítja a felfedési folyamatot, és az AI az innováció élvonalába tartozik.³⁰ Az ügyfelek ma már gyakran kérik fel a peres ügyvédeket, hogy állítsanak fel e-discovery relevancia-hipotéziseket, és alkalmazzanak prediktív kódolási modelleket (a TAR egy fajtája) az elektronikus információk feltárásához.³¹ Ebben a folyamatban a peres ügyészek először meghatározzák a keresendő kulcsszavakat, és azonosítják az áttekintendő dokumentumok kezdeti csoportját.³² Ezután a dokumentumokat vizsgáló ügyvédek felülvizsgálják, kódolják és pontozzák a dokumentumok kezdeti halmazát bizonyos kulcsszavak előfordulása alapján a dokumentum relevanciájával kapcsolatban.³³ A felülvizsgálat során az e-discovery ügyvédek felügyelt tanulási algoritmusokat képeznek és modelleznek a dokumentumok osztályozására, amelyek a dokumentumokat felülvizsgáló ügyvédek döntései alapján osztályozzák a dokumentumokat a kezdeti dokumentumhalmazban.³⁴ Más szóval, az algoritmus a valódi ügyvédek döntéseinek elemzésével és megismétlésével tanulja meg, hogy mely dokumentumok relevánsak.³⁵ Ezen túlmenően a prediktív kódolási modellek több millió felderíthető dokumentumot képesek relevancia alapján osztályozni.³⁶

21 Lásd NILSSON, *Supra* note 19, 388-89. o.

22 Lásd NICK BOSTROM, SUPERINTELLIGENCE: (232017).

23 TEGMARK, *fenti* megjegyzés 5.

24 NILSSON, *fenti* megjegyzés 19.

25 *Id.*

26 *Az OpenAI-ról*, OPENAI <https://openai.com/about/> (utolsó látogatás 2019.10., május).

27 Lásd általában Nick Bostrom, *Are You Living in A Computer Simulation?*, FILOZÓFIA 53Q. (2003)211.,243

28 ALLEN & CHAN, *Supra* note 12; lásd még HEMANT TANEJA, UNSCALED: HOW AI AND NEW GENERATION OF UPSTARTS ARE CREATING THE ECONOMY OF THE FUTURE (12018).

29 RICHARD SUSSKIND, TOMORROW'S LAWYERS (112d ed. 2017).

30 Scott D. Cessar, Christopher R. Opalinski, & Brian E. Calla, *Controlling Electronic Discovery Costs: Cutting "Big Data" Down to Size*, ECKERT SEAMANS (2013. március 5.), <https://www.eckertseamans.com/publications/controlling-electronic-discovery-costs-cutting-big-data-down-to-size/>; lásd még Nicholas Barry, *Man Versus Machine Review: The Showdown Between Hordes of Discovery Lawyers and a Computer-Utilizing Predictive-Coding Technology*, VAND15. J. ENT. & TECH. L. (2013)343.,344

31 KEVIN D. ASHLEY, ARTIFICIAL INTELLIGENCE AND LEGAL ANALYTICS 240-42 (2017).

32 Barry, *Supra* note at 30,351.

33 GORDON V. CORMACK & MAURA R. GROSSMAN, EVALUATION OF MACHINE-LEARNING PROTOCOLS FOR TECHNOLOGY-ASSISTED REVIEW IN ELECTRONIC DISCOVERY (1542014), <http://plg2.cs.uwaterloo.ca/~gvcormack/calstudy/study/sigir2014-cormackgrossman.pdf>.

34 Barry, *Supra* note at 30,354.

35 *Id.*

36 *Lásd pl.* ASHLEY, *Supra* note at 31,250.

A másik példa a mesterséges intelligencia miatt gyorsan fejlődő iparágra az egészségügy.³⁷ Egy évtized múlva az egészségügyi ágazat a mesterséges intelligenciának köszönhetően egészen másképp fog kinézni, mint ma.³⁸ Jelenleg a nagy adatok által vezérelt mesterséges intelligencia érezhető elmozdulást hoz létre az orvosi gyakorlatban a tömeges ellátástól a személyre szabott ellátás felé.³⁹ A modern kórházakban praktizáló egészségügyi szakemberek ugyanis ma már elektronikus adatbázisokban, elektronikus egészségügyi nyilvántartásokban ("EHR") tárolják a betegek adatait.⁴⁰ Ez lehetővé teszi, hogy a gépi tanuló algoritmusok elemezzék a betegek egészségügyi adatait, és drasztikusan javítsák a betegellátást.⁴¹ Ezek az adatvezérelt erőforrások nemcsak azt teszik lehetővé, hogy az orvos gyakorlatilag mindent tudjon a beteg kórtörténetéről anélkül, hogy valaha is találkozna a beteggel, hanem az orvosi munka segítségével drasztikusan csökkentik az egészségügyi ellátással kapcsolatos költségeket is.⁴² Például 2016-ban a Stanford kutatói olyan mesterséges intelligenciát fejlesztettek ki, amely pontosabban tudta diagnosztizálni a tüdőrákot, mint az emberi patológusok.⁴³ Egy másik példa a D-Wave adiabatikus kvantumszámítógépe, amely képes a rákdiagnosztikához szükséges gépi tanulási algoritmusok futtatására.⁴⁴ Röviden, az EHR-ek, a nagyméretű adatok és a mesterséges intelligencia átalakítják az egészségügyi ellátást.⁴⁵

A mesterséges intelligencia okozta zavarok harmadik példája a védelmi iparban zajlik. A mesterséges intelligencia már most is alapvető eszköz a kiberbiztonságban.⁴⁶ Mike Rogers admirális, a Nemzetbiztonsági Hivatal igazgatója szerint a mesterséges intelligencia és a gépi tanulás a kiberbiztonság jövőjének alapját képezi.⁴⁷ Márciusban 2,2017, jelentést tettek közzé a Fehér Házban, amely szerint orosz programozók mesterséges intelligenciával végrehajtott kibertámadást indítottak a Védelmi Minisztérium több mint 10 000 alkalmazottjának személyes közösségi média fiókjai ellen.⁴⁸ Emellett a modern hadviselésben a mesterséges intelligenciát a csatatéren is használják.⁴⁹ Például a haditengerészeti hajókra telepített amerikai Phalanx rakétavédelmi rendszer mesterséges intelligenciát használ az ellenséges rakéták és repülőgépek által jelentett fenyegetések észlelésére, követésére és támadására.⁵⁰ A kereskedelmi AI-rendszerekkel való terrorista visszaélés azonban komoly problémát jelent.⁵¹ Terrorista szervezetek már most is használják a drónokban lévő mesterséges intelligencia-rendszereket robbanóanyagok szállítására és balesetek okozására.⁵² A szűk körű mesterséges intelligencia továbbra is megváltoztatja az olyan professzionális iparágak működését, mint a jog, az egészségügy és a védelem.⁵³ Számos AI-kutató hivatkozott megfigyelhető

37 TEGMARK, *Supra* note at 5,102.

38 TANEJA, *Supra* note at 28,73.

39 *Id.*

40 Kate Monica, *Apple EHR Patient Data Viewer Now in Use at Health39 Systems*, EHRINTELLIGENCE (2018. április 2.), <https://ehrintelligence.com/news/apple-ehr-patient-data-viewer-now-in-use-at-39-health-sys-tems>.

41 *Lásd* XIAOQIAN JIANG ET AL., A PATIENT-DRIVEN ADAPTIVE PREDICATION TECHNIQUE TO IMPROVE PERSONALIZED RISK ESTIMATION FOR CLINICAL DECISION SUPPORT 137 (2012).

42 *Lásd* Alvin Rajkomar et al., *Scalable and Accurate Deep Learning with Electronic Health Records*, NATURE PARTNER J. (2018), <https://www.nature.com/articles/s41746-018-0029-1.pdf>.

43 *Lásd* Lloyd Minor, *Crunching the Image Data Using Artificial Intelligence to Look at Biopsies*, STANFORD MED. (2017), <https://stanmed.stanford.edu/2017summer/artificial-intelligence-could-help-diagnose-can-cer-predict-survival.html>.

44 *Lásd* Brian S. Haney, *Quantum Machine Learning Cancer Diagnostics*, GITHUB (2019. február 24.), https://github.com/Bhaney44/Leap/blob/master/Quantum_Machine_Learning_Cancer_Diagnostics.py.

45 *Id.*

46 *Lásd általánosságban* BRUNDAGE ET AL. „*supra* note 5.

47 ALLEN & CHAN, *Supra* note at 12,18.

48 *Lásd* Massimo Calbresi, *Inside Russia's Social Media War on America*, TIME (2017. május), <http://time.com/4783932/inside-russia-social-media-war-america/>.

49 *Lásd* *United States Navy Fact File: MK15 - Phalanx Close-In Weapons System (CIWS)*, U.S. DEP'T

NAVY (utolsó látogatás 201913., május),

http://www.navy.mil/navydata/fact_display.asp?cid=2100&tid=487&ct=2.

50 TEGMARK, *Supra* note at (5,1112017).

51 *Lásd általánosságban* BRUNDAGE ET AL., *supra* note 5.

52 *Lásd id.*

53 *Lásd általában* TEGMARK, *fenti megjegyzés. 5.*

az információtechnológiai ár- és teljesítménykinetika történelmi mintázatai, amelyek alátámasztják azt az érvet, hogy a mesterséges intelligencia technológiák fejlődési üteme a vártnál sokkal gyorsabban fog végbemenni.⁵⁴ Sőt, ezek a kutatók feltételezik, hogy a mesterséges intelligencia technológiák továbbra is gyorsuló ütemben fognak fejlődni.⁵⁵

B. A GYORSULÓ MEGTÉRÜLÉS TÖRVÉNYE

A gyorsuló megtérülés törvénye ("LOAR") kimondja, hogy az információs technológia alapvető mérőszámai általában kiszámítható és exponenciális pályát követnek.⁵⁶ Az információs technológiák valóban exponenciális módon épülnek egymásra; ezt a jelenséget Moore törvényének nevezték el, és könnyen mérhető a legtöbb olyan folyamatban, ahol az információ mintái fejlődnek.⁵⁷ A LOAR a számítástechnika⁵⁸ ára és teljesítményére való alkalmazását írja le, és Gordon Moore, az Intel alapítója javasolta 1965-ben.⁵⁹ A Moore-törvény azt jósolja, hogy másfél évente a számítógépek feldolgozási teljesítménye megduplázódik, miközben a költségek a felére csökkennek.⁶⁰ Általánosságban azt jelenti, hogy az informatika teljesítménye másfél évente megduplázódik.⁶¹ Az elmúlt ötvenhárom év igazolta Gordon Moore jóslatát; ⁶²egy okostelefon ma nagyobb számítási teljesítményt nyújt, mint az egész NASA 1969-ben - amikor az Apollo-11 leszállt a Holdra.⁶³ A Moore-törvényt a mesterséges intelligenciára alkalmazva sok mesterséges intelligencia kutató úgy véli, hogy jelenleg a szuperintelligens mesterséges intelligencia kifejlesztésének küszöbén állunk.⁶⁴

Irving J. Good 1965-ben mutatta be először a szuperintelligencia fogalmát.⁶⁵ Good kijelentette: "[I]etezzük az ultraintelligens gépet úgy, mint olyan gépet, amely messze felülmúlja bármely ember minden szellemi tevékenységét, legyen az bármilyen okos is".⁶⁶ Good szerint "[a]lévén, hogy a gép tervezése az egyik ilyen intellektuális tevékenység, egy ultraintelligens gép még jobb gépeket tudna tervezni; ekkor kétségtelenül intelligencia-robbanás következne be, és az ember intelligenciája messze lemaradna".⁶⁷ Good valóban azt jósolta, hogy az első ultraintelligens gép lesz "az utolsó találmány, amelyet az embernek valaha is meg kell tennie".⁶⁸ A közelmúlt tudósai átvették Good elemzését, és hasonlóan definiálták a szuperintelligenciát. Nick Bostrom oxfordi professzor például úgy definiálja a szuperintelligenciát, mint "minden olyan értelmet, amely gyakorlatilag minden érdeklődési területen jelentősen meghaladja az emberek kognitív teljesítményét".⁶⁹ Max Tegmark szerint a szuperintelligencia "[g]enerális intelligencia, amely messze meghaladja az emberi szintet".⁷⁰

54 Lásd BOSTROM, *Supra* 22. lábjegyzet, a következő címen 85.

55 Lásd RAY KURZWEIL, HOGYAN KÉSZÜLJÜK AZ ELMET (2502012).

56 Lásd *id.*

57 Lásd *id.* 256.

58 Lásd MARTINE ROTHBLATT, VIRTUÁLIS EMBER (482014).

59 Lásd KURZWEIL, *Supra* note at 55,251.

60 Lásd SUSSKIND, *Supra* note at 29,11.

61 Lásd ROTHBLATT, *Supra* note at 58,28.

62 Lásd KURZWEIL, *Supra* note at 55,251.

63 Lásd MICHIO KAKU, A JÖVŐ FIZIKÁJA (232011).

64 Lásd általában BOSTROM, 22. lábjegyzet; lásd még KURZWEIL, 55. lábjegyzet; lásd még TEGMARK, 55. lábjegyzet. megjegyzés 5.

65 Lásd általában Irving J. Good, *Speculations Concerning the First Ultraintelligent Machine*, AD-6 VANCES IN COMPUTERS (311966).

66 *Id.* 33.

67 *Id.*

68 *Id.*

69 BOSTROM, *Supra* note at 22,22.

70 Lásd TEGMARK, *fenti* jegyzet, lásd:5, TEGMARK, *fenti* jegyzet. 39.

A LOAR alkalmazása a mesterséges intelligenciára azt bizonyítja, hogy a szűk értelemben vett mesterséges intelligenciáról az AGI-ra és a szuperintelligenciára való áttérés sokkal közelebb lehet, mint azt általában gondolják.⁷¹ Egyelőre az AGI legkorábbi becslések szerint 2029-re várható.⁷² Ray Kurzweil szerint a XXI. század a LOAR-nak köszönhetően olyan technológiai fejlődést és innovációt hoz, amely ma még éveknek 20,000 tűnhet.⁷³ Emellett Bostrom és Eliezer Yudkowsky, a mesterséges intelligencia teoretikusa azt jósolta, hogy a közvélemény a mesterséges intelligencia antropomorfizmusa miatt a mesterséges intelligencia fejlődésének gyors kinetikáját fogja érzékelni.⁷⁴ A mesterséges intelligencia antropomorfizmusa az emberi intelligenciaszintek nem emberi entitásoknak való tulajdonítására utal.⁷⁵ Az emberek az intelligenciaspektrum szélsőséges végpontjainak tekinthetik a falusi bolondot és Albert Einsteint,⁷⁶ de a kettő közötti különbség nagyobb relatív skálán valójában *csekély*.⁷⁷ Így egy mesterséges intelligencia-rendszer fejlődése a falusi idióta intelligenciájától Einstein intelligenciájáig, az AGI intelligenciájáig és végül a szuperintelligenciáig gyorsabb lehet a vártnál.⁷⁸

Érdekes módon ezeket az előrejelzéseket alátámasztja az a hatalmas mennyiségű információ, amelyet az emberek a digitális korszak hajnalán kezdtek gyűjteni.⁷⁹ Valójában az emberek által gyűjtött információk mennyisége is egyre gyorsul.⁸⁰ Az adatok, amelyeket a világról szóló információk digitális reprezentációjaként határoznak meg,⁸¹ elképesztő sebességgel keletkeznek. Kétnaponta az emberek több mint öt kvintillió bajtnyi adatot hoznak létre, vagyis annyi adatot, mint a civilizáció hajnalától kezdve egészen a 2003.⁸²

Michael Kremer harvardi professzor és közgazdász szerint "az emberi fejlődés alapvető mozgatórugója nem a nyersanyagok, hanem a problémák technológiai megoldása".⁸³ A mesterséges intelligencia kontextusában az adatok a technológiai fejlődés hajtóereje az emberi programozók helyett.⁸⁴ A technológiai megoldások hajtóereje pedig az a felismerés, hogy minden információ számokként ábrázolható.⁸⁵ Az egy adott problémához rendelkezésre álló adatok mennyisége és típusa nagyban meghatározza a fejleszhető AI-rendszerek erejét.⁸⁶ A LOAR tehát mélyreható hatással lesz az AI fejlődésére az AGI és a szuperintelligencia irányába.

Egyesek azonban azt állítják, hogy az AGI talán soha nem fog megvalósulni.⁸⁷ A Microsoft néhai társalapítója, Paul Allen például azt állítja, hogy a tudományos fejlődés szabálytalan, és feltételezi, hogy a XXI. század végére az embereknek még

71 Lásd *id.* 157.

72 Lásd KURZWEIL, *Supra* note at 55,261.

73 Lásd Ray Kurzweil, *The Law of Accelerating Returns*, in KURZWEIL NETWORK (2001), <http://www.kurzweil.net/the-law-of-accelerating-returns>.

74 Lásd BOSTROM, *Supra* 22. lábjegyzet, a következő címen 85.

75 Lásd Eliezer Yudkowsky, *A mesterséges intelligencia mint a globális kockázat pozitív és negatív tényezője*, MACHINE INTELLIGENCE RES. INST., (212008), <https://intelligence.org/files/AIPosNegFactor.pdf>.

76 Lásd BOSTROM, *Supra* 22. lábjegyzet, a következő címen 85.

77 Lásd Yudkowsky, *Supra* note. 75.

78 Lásd BOSTROM, *Supra* 22. lábjegyzet, a következő címen 86.

79 Lásd ETHEM ALPAYDIN, MACHINE LEARNING (112016).

80 Lásd SUSSKIND, *Supra* note at 2911.

81 ALPAYDIN, *Supra* note at 79,3.

82 SUSSKIND, *Supra* note at 29,11.

83 Michael Kremer, "Népességnövekedés és technológiai változás: Kr.e. egymillió évtől 1990 a.," 108Q. J. OF ECON. 3 (1993). (idézi SAIFEDEAN AMMOUS, THE BITCOIN STANDARD: THE DECENTRALIZED ALTERNATIVE TO CENTRAL BANKING (2018).

84 ALPAYDIN, *Supra* note at 79,12.

85 *Id.*

86 SEBASTIAN RASCHKA & VAHID MIRJALILI, PYTHON MACHINE LEARNING, (22d. szerk. 2017).

87 Paul G. Allen & Mark Greaves, *The Singularity Isn't Near*, MIT TECH. REV. (201112., október), <https://www.technologyreview.com/s/425733/paul-allen-the-singularity-isnt-near/>.

AGI elérése.⁸⁸ Másrészt Max Tegmark szerint a vita alapvető igazsága - hogy az emberiség valaha is létrehozza-e az AGI-t - továbbra is bizonytalan.⁸⁹ Tegmark azonban azt is kifejti, hogy a legtöbb AI-szakértő szerint az AGI 2047 körül fog megvalósulni.⁹⁰ Egy tudós szerint az AI hatásával kapcsolatos kérdések csak még sürgetőbbé válnak, ahogy közeledünk az exponenciális inflexiós ponthoz, és a növekedés hirtelen és drámai függőleges pályára áll.⁹¹ Egyelőre az a kérdés, hogy a társadalom közeledik-e ehhez a fordulóponthoz, vagy még mindig a lassabb, fokozatos fejlődési szakaszban van.⁹² Ma az emberiség legegységesebb útja az AGI létrehozása felé a mély megerősítő tanulás.

II. AGI FEJLESZTÉS

A gépi tanulás a mesterséges intelligencia egyik részterülete, amely a gépek azon képességére összpontosít, hogy megtanulják és lemásolják az emberi elmével kapcsolatos kognitív viselkedést.⁹³ A gépi tanulás általában adatbányászatot, mintafelismerést és természetes nyelvi feldolgozást foglal magában.⁹⁴ Ezek a technikák az elmúlt években egyre népszerűbbé váltak, mivel az internet megjelenése óta az emberek által előállított és összegyűjtött adatok mennyisége robbanásszerűen megnőtt.⁹⁵ A gépi tanulás legújabb áttörése a mély megerősítéses tanulás.⁹⁶ A mély megerősítéses tanulás a gépi tanulás két hagyományos modelljét - a felügyelt tanulást és a megerősítéses tanulást - ötvözi, hogy az algoritmusok az embertől függetlenül tanulhassanak.⁹⁷

A mesterséges intelligencia szabályozásával kapcsolatos legtöbb tanulmány a gépi tanulás felügyelt vagy nem felügyelt módszereire összpontosít, mivel 2014-ig ez volt a gépi tanulás egyetlen két típusa, amelyet széles körben használtak.⁹⁸ Valójában a legtöbb jogi tanulmány a mély neurális hálózatokra, a felügyelt tanulási algoritmusok egyik típusára összpontosít.⁹⁹ A Google azonban 2013-ban kifejlesztette a tanulás egy új típusát, a "mély megerősítéses tanulást", amelyet később szabadalmaztatott.¹⁰⁰ Az 1980-as években úttörő szerepet játszó megerősítéses tanulás egy olyan gépi tanulási technika, amelyet a behaviorista pszichológia ihletett, ahol egy intelligens ágens hajlamát, hogy bizonyos módon cselekedjen, egy jutalmazási struktúra befolyásolja.¹⁰¹ Az intelligens ágens olyan entitás, amely szenzorok segítségével információkat gyűjt a környezetéről, majd ezeket az információkat feldolgozva dönt arról, hogyan reagáljon a környezetére.¹⁰² A mély megerősítéses tanulás a megerősítést kombinálja

88 *Id.*

89 TEGMARK, 5. lábjegyzet, 54.

pont. 90 *Id.* 157.

91 Michael Guihot et al., *Nudging Robots: Innovatív megoldások a mesterséges intelligencia szabályozására*, VAND20. J. ENT. & TECH. L. (2017)385.,400

92 *Lásd id.*

93 *Lásd általánosságban* ALPAYDIN, *fenti* megjegyzés. 79.

94 *Lásd* Michael Simon et al., *Lola v. Skadden and the Automation of the Legal Profession*, 20 YALE J.L. & TECH (234,2018253) (*idézi* Bernard Marr, *What Everyone Should Know About Cognitive Computing*, FORBES (2016. március 23., 3:28), <https://www.forbes.com/sites/bernardmarr/2016/03/23/what-everyone-should-know-about-cognitive-computing/#5630f9005088>).

95 *Lásd id.* 252.

96 *Lásd* TEGMARK, *fenti* jegyzet, lásd:5, TEGMARK, *fenti* jegyzet. 85.

97 *Lásd id.*

98 *Lásd id.* 83.

99 *Lásd általában* Calo, 6. lábjegyzet; *lásd még* John O. McGinnis, *Accelerating AI*, NW104. U. L. REV. 1253 (2010).

100 '862 Application, *fenti* megjegyzés 15.

101 RICHARD S. SUTTON & ANDREW G. BARTO, REINFORCEMENT LEARNING: AN INTRODUCTION (552017); *lásd még* TEGMARK, *supra* note at5, (2017). 85.

102 TEGMARK, *Supra* note at 5,84.

tanulás mély neurális hálózatok használatával.¹⁰³ A mély megerősítéses tanulás olyan megerősítéses tanulási algoritmusra utal, amely mély neurális hálózatot használ függvényközelítőként, és amelyet e rész későbbi részében ismertetünk.¹⁰⁴ Ez a rész először a mély neurális hálózatokat ismerteti. Másodszor, ez a rész a megerősítő tanulást ismerteti. Harmadszor, ez a rész elmagyarázza a mély megerősítő tanulást.

A. MÉLY NEURÁLIS HÁLÓZATOK

Az emberi agy "neuronoknak" nevezett feldolgozó egységekből áll.¹⁰⁵ Az agy minden egyes neuronja más neuronokkal szinapszisoknak nevezett struktúrákon keresztül kapcsolódik.¹⁰⁶ Egy biológiai neuron dendritekből áll - más neuronoktól érkező különböző elektromos impulzusok vevői -, amelyek a neuron sejttestében gyűlnek össze.¹⁰⁷ Amint a neuron sejtteste elegendő elektromos energiát gyűjtött össze, hogy meghaladjon egy küszöbértéket, a neuron elektromos töltést továbbít az agy más neuronjainak a szinapszisokon keresztül.¹⁰⁸ A biológiai agyban ez az információátvitel adja az alapját a modern ideghálózatok működésének.¹⁰⁹

A mesterséges neuronok lényegében a biológiai neuronról mintázott logikai kapuk.¹¹⁰ Mind a mesterséges, mind a biológiai neuronok különböző forrásokból kapnak bemenetet, és a bemeneti információt egyetlen kimeneti értékre képezik le.¹¹¹ A mesterséges neuronhálózat egymással összekapcsolt mesterséges neuronok csoportja, amelyek képesek egymás viselkedését befolyásolni.¹¹² A mesterséges neurális hálózatban a neuronok a biológiai agyban lévő szinapszisok erősségét modellező súlykoefficiensekkel vannak összekapcsolva.¹¹³ A neurális hálózatokat nagy adathalmazok felhasználásával képzik.¹¹⁴ A képzési folyamat lehetővé teszi a súlykoefficiensek beállítását, hogy a neurális hálózat kimenete vagy előrejelzése pontos legyen.¹¹⁵ Miután a neurális hálózatot betanították, új adatokat táplálnak a hálózaton keresztül, hogy előrejelzéseket készítsenek.¹¹⁶

1957-ben Frank Rosenblatt közzétett egy algoritmust - a perceptront -, amely automatikusan megtanulja az optimális súlykoefficienseket egy mesterséges neurális hálózat számára.¹¹⁷ A perceptron modellt az alábbiakban szemléltetjük.¹¹⁸

103 Fei-Fei Li, Justin Johnson, & Serena Yeung, *14. előadás: Mély megerősítéses tanulás*, STANFORD U. SCH. OF ENG'G (2017), <https://www.youtube.com/watch?v=lvoHnicueoE> (utolsó hozzáférés: 201913., május).

104 *Id.*

105 ALPAYDIN, *Supra* note at 79,86.

106 *Id.*

107 RASCHKA & MIRJALILI, *uo.* 86. lábjegyzet, at 18.

108 *Id.*

109 ALPAYDIN, *Supra* note at 79,86.

110 KURZWEIL, *Supra* note at 55,38.

111 JOHN D. KELLEHER & BRENDAN TIERNEY, ADATTUDOMÁNY (1312018).

112 TEGMARK, *Supra* note at 5,72.

113 ALPAYDIN, *Supra* note at 79,88.

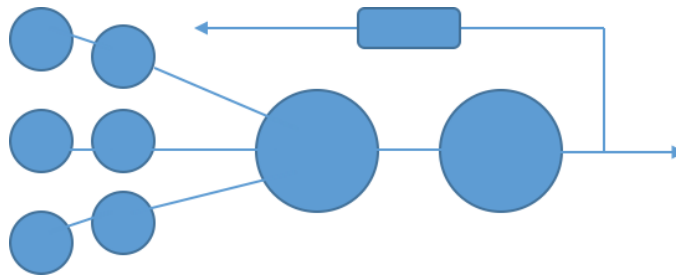
114 *Id.* 89.

115 *Id.*

116 KELLEHER, *Supra* note 111, at 127.

117 RASCHKA & MIRJALILI, *uo.* 86. lábjegyzet, at 18.

118 *Id.*



A perceptronban a bal szélső három kör a bemeneti értékeket $x_j \dots m$, a hozzá tartozó súlyértékek pedig $w_{j \dots m}$ bemeneti értékektől jobbra lévő három kör.¹¹⁹ A bemeneti értékek és a súlyértékek aggregálásra kerülnek, jellemzően egy összegző egyenlet segítségével, amelyet az első nagy kör (balról jobbra) képvisel.¹²⁰ A második nagy kör a küszöbfüggvényt jelképezi, egy előre meghatározott értéket, amelynek túllépése esetén a modell 1 kimenetet jelez.¹²¹ Ha a küszöbfüggvényt nem lépi túl, a modell 0-t ad ki.¹²² A kimenetet a jobbra mutató nyíl jelzi.¹²³ A modell tetején lévő doboz a hibafüggvényt ábrázolja.¹²⁴ Abban az esetben, ha a modell kimenete hibás, akkor a hibafüggvény lép működésbe.¹²⁵ Ha a hibafüggvény működésbe lép, a súlyértékek a perceptron tanulási szabály szerint frissülnek.¹²⁶ A perceptron tanulási szabály formális ábrázolása a következő

mint: $\Delta w_j = \eta (y^{(i)} - \hat{y}^{(i)}) x_j$, ahol η a tanulási ráta, $y^{(i)}$ a valós osztálycímke a az i -edik gyakorló minta, és $\hat{y}^{(i)}$ a megjósolt osztálycímke.¹²⁷ A valódi osztálycímke a kimeneti címke, a megjósolt osztálycímke pedig a perceptron kimenete.¹²⁸

Minden neurális hálózatnak van egy bemeneti és egy kimeneti rétege.¹²⁹ A bemeneti és kimeneti réteg között azonban a neurális hálózatok több rejtett réteget tartalmaznak.¹³⁰ A rejtett rétegek száma változhat, és az adott modelltől függ.¹³¹ Fontos megjegyezni, hogy míg a perceptron modellek általában lineáris osztályozási feladatokra korlátozódnak, ez a korlátozás nem vonatkozik a többrétegű hálózatokra.¹³² A többrétegű perceptron modell ugyanis egy univerzális approximátor, azaz olyan algoritmus, amely elegendő neuron esetén bármilyen függvényt képes a kívánt pontossággal közelíteni.¹³³ A mély neurális hálózat olyan hálózat, amely több rejtett réteggel rendelkezik.¹³⁴ Ez lehetővé teszi, hogy a neurális hálózat több absztrakciós réteget vegyen figyelembe.¹³⁵ Az alábbi ábra egy mély neurális hálózat egyszerű modellje.¹³⁶

119 *Id.* 19.

120 *Lásd* ALPAYDIN, *fenti* jegyzet, lásd:79, ALPAYDIN, *fenti* jegyzet. 89.

121 *Lásd Id.*

122 *Lásd* RASCHKA & MIRJALILI, 86. lábjegyzet, a következő címen 18.

123 *Lásd* KURZWEIL, *Supra* note at 55,132.

124 *Lásd* RASCHKA & MIRJALILI, 86. lábjegyzet, a következő címen 18.

125 *Lásd* ALPAYDIN, *fenti* jegyzet, lásd:79, ALPAYDIN, *fenti* jegyzet. 90.

126 *Lásd id.*

127 *Lásd* RASCHKA & MIRJALILI, 86. lábjegyzet, a következő címen 21.

128 *Lásd id.* 22.

129 *Lásd* KURZWEIL, *Supra* note at 55,132.

130 *Lásd* ALPAYDIN, *fenti* jegyzet, lásd:79, ALPAYDIN, *fenti* jegyzet. 100.

131 *Lásd* KELLEHER & TIERNEY, 111,132. lábjegyzet.

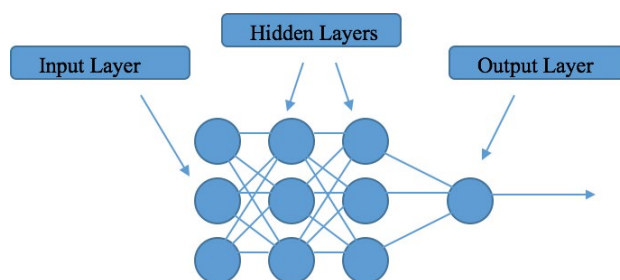
132 *Lásd* ALPAYDIN, *fenti* jegyzet, lásd:79, ALPAYDIN, *fenti* jegyzet. 99.

133 *Lásd id.*

134 *Lásd* TEGMARK, *fenti* jegyzet, lásd:5, TEGMARK, *fenti* jegyzet. 76.

135 ALPAYDIN, *Supra* note at 79,88.

136 KELLEHER & TIERNEY, *Supra* note at 11, (a következő idézetben szereplő illusztráción alapuló 132modell).



Minden neuron egy rejtett egységet képvisel egy rétegben, és a modell egy komplex tulajdonságát határozza meg.¹³⁷ A rejtett egységek megfelelnek a megfigyelt, de közvetlenül nem megfigyelt rejtett tulajdonságoknak.¹³⁸ A rejtett egységek egymást követő rétegei pedig a jellemzők absztrakciójának növekvő rétegeinek felelnek meg.¹³⁹

Valójában a rejtett egységek minden egyes rétege egy-egy jellemző-kivonó funkcióként működik, mivel kissé bonyolultabb jellemzők elemzését biztosítja.¹⁴⁰ A jellemző-kivonás a dimenziócsökkentés módszere - a bemeneti jellemzők csökkentésének módszere -, amely lehetővé teszi a nyers bemenet átalakítását kimenetű oly módon, hogy az adattudósok megfigyelhessék az adatok rejtett jellemzőit.¹⁴¹ A későbbi rejtett egységek a rejtett jellemzőket a korábbi jellemzők kombinálásával vonják ki a bemeneti tér egy kissé nagyobb részében.¹⁴² A kimeneti réteg a teljes bemenetet figyeli, hogy végső előrejelzést készítsen.¹⁴³ Más szóval, a mély neurális hálózatok a kezdeti bemenetük bonyolultabb függvényeit tanulják meg, amikor minden rejtett réteg kombinálja az előző réteg értékeit.¹⁴⁴ Emellett a mély neurális hálózatok számos kontextusban kiválóan alkalmasnak bizonyultak előrejelzések készítésére.¹⁴⁵ Ezek a modellek azonban a tanuláshoz adatokra és legalább minimális emberi beavatkozásra van szükség a tanulási folyamat felügyeletéhez.¹⁴⁶ A megerősítéses tanulás egy újabb gépi tanulási technika, amely egyiket sem igényli.¹⁴⁷

B. MÉLY MEGERŐSÍTÉSES TANULÁS

A megerősítéses tanulás egyfajta gépi tanulási technika, amelyet a behaviorista pszichológia ihletett.¹⁴⁸ Formálisan a megerősítéses tanulás egy ágens-környezet kölcsönhatáson keresztül írható le, a Markov-döntési folyamattal ("MDP").¹⁴⁹ Az alábbi modell az ágens-környezet kölcsönhatást írja le egy MDP-ben.¹⁵⁰

¹³⁷ ALPAYDIN, *Supra* note at 79,100.

¹³⁸ *Id.*

¹³⁹ *Id.*

¹⁴⁰ KELLEHER & TIERNEY, *Supra* note at 111,135.

¹⁴¹ ALPAYDIN, *Supra* note at 79,102.

¹⁴² *Id.*

¹⁴³ KURZWEIL, *Supra* note at 55,132.

¹⁴⁴ ALPAYDIN, *Supra* note at 79,104.

¹⁴⁵ *Lásd általában* ASHLEY, *fenti* megjegyzés. 31.

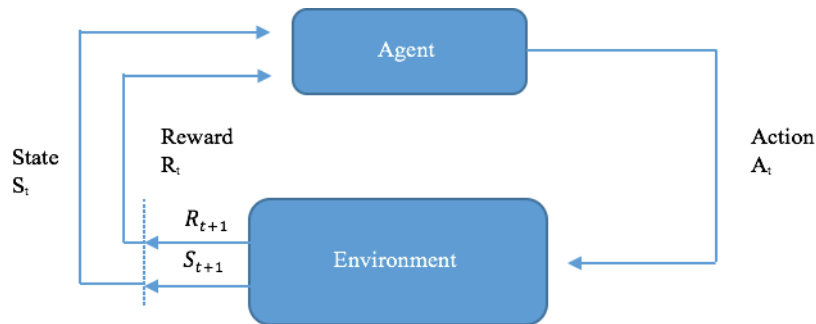
¹⁴⁶ ALPAYDIN, *Supra* note at 79,106.

¹⁴⁷ *Id.* 127.

¹⁴⁸ SUTTON & BARTO, *Supra* note at 101,38.

¹⁴⁹ Alex Kendall et al., *Learning to Drive in a Day*, (1201811, szeptember) (nem publikált tanulmány) (elérés: <https://arxiv.org/abs/1807.00412>).

¹⁵⁰ SUTTON & BARTO, *Supra* note 101, 38. o. (az előző idézetben szereplő illusztráción alapuló modell).



A környezet minden egyes időpontban, amikor a környezet létezik, állapotokból áll.¹⁵¹ Az ágens cselekedetei az egyes állapotokban meghatározzák a környezet valószínűségi fejlődését.¹⁵²

Kezdetben az ágens elé kerül a környezet állapota, amely több lehetséges cselekvést is tartalmaz.¹⁵³ Ezután az ágens végrehajt egy akciót az adott állapotban, és továbblép a környezet következő állapotába, ahol jutalmat kap.¹⁵⁴ Az ágens az ágens politikája alapján választja ki, hogy melyik akciót hajtja végre, amikor egy állapotot mutatunk neki.¹⁵⁵ A politika az a mód, ahogyan az ágens döntéseket hoz vagy cselekvéseket választ egy állapoton belül.¹⁵⁶ Például egy nagyfokú integritással rendelkező személynek van egy politikája, amely rutinszerűen irányítja a döntéshozatalát, hogy a helyes dolgot válassza, amikor etikai dilemmákkal szembesül. Hasonlóképpen, egy kapzsi személynek olyan politikája van, amely rutinszerűen irányítja a döntéshozatalát, hogy a legmagasabb dollárértéket hozó cselekvést válassza. A politika célja, hogy az ágens úgy haladjon előre a környezetben, hogy maximalizálja a jutalmat.¹⁵⁷

Egy érték-funkció határozza meg az s állapotban való tartózkodás és a politika követésének értékét.

π a környezet végső állapotáig, amelyet végállapotnak nevezünk.¹⁵⁸ A végállapot zárja le az epizódot, amely egy környezet összes állapotából áll.¹⁵⁹ Az π politika π végrehajtásának várható értékét az s állapot mellett $V^\pi(s)$ jelöli.¹⁶⁰ Egy MDP kontextusában az V^π értékfüggvény egyenlő az π politika végrehajtásáért járó diszkontált jutalmak várható összegével:¹⁶¹

$$V^\pi(s) = E[R(s_0) + \gamma R(s_1) + \dots | s_0 = s, \pi(s)]$$

A várható jövőbeni jutalmakat γ diszkontfaktorial diszkontáljuk.¹⁶² A diszkontfaktor jellemzően a következő: $0 < \gamma < 1$.¹⁶³ Ez lehetővé teszi az értékfüggvény meghatározását

151 Li, Johnson, & Yeung, *Supra* note. 103.

152 MYKEL J. KOCHENDERFER, DÖNTÉSHOZATAL BIZONYTALANSÁG MELLETT (772015).

153 SUTTON & BARTO, *fenti* 101,39. lábjegyzet.

154 KOCHENDERFER, *Supra* note, 152,77. o.; *lásd még* SUTTON & BARTO, *Supra* note, 101, 77. o. 39.

155 SUTTON & BARTO, *fenti* 101,39. lábjegyzet.

156 KOCHENDERFER, *supra* note at 152,77.

157 SUTTON & BARTO, *fenti* 101,50. lábjegyzet.

158 Li, Johnson, & Yeung, *Supra* note. 103.

159 *Id.*

160 KOCHENDERFER, *supra* note at 152,80.

161 Ahmad El Sallab et al., *Deep Reinforcement Learning Framework for Autonomous Driving*, (Apr. 20178.), (nem publikált tanulmány) (elérhető a <https://arxiv.org/pdf/1704.02532.pdf> oldalon. *Lásd az A. függelékben a jelölések összefoglalóját.*)

162 SUTTON & BARTO, 101,92. lábjegyzet.

163 *Id.*

véges időben, és lehetővé teszi, hogy a jelenbeli jutalmak értéke értékesebb legyen, mint a jövőbeli jutalmaké.¹⁶⁴ Az $\pi^*(s)$ optimális politikát úgy definiáljuk, hogy az a politika, amely maximalizálja a várható értéket a többi politikához képest.¹⁶⁵ Az MDP modell célja az optimális politika megtalálása:¹⁶⁶

$$\pi^*(s) = \underset{\pi}{\operatorname{arg\,max}} V^\pi(s)$$

Egy adott MDP optimális politikájának megtalálását általában Q-tanulással oldják meg.¹⁶⁷ A Q-tanulás ezt a problémát egy Q-értékfüggvény maximalizálásával oldja meg: $Q(s, a)$.¹⁶⁸ A Q-értékfüggvény egy állapot-cselekvés pár értékét írja le.¹⁶⁹ Valójában a Q-tanuló algoritmus célja az optimális Q-értékfüggvény felfedezése.

Q^* bármely állapot-cselekvés párosra.¹⁷⁰ A Bellman-egyenlet kifejezi egy állapot értéke és az azt követő állapotok értékei közötti kapcsolatot.¹⁷¹ Az algoritmus a Q-értékfüggvény konvergenciájáig örökké folytatódik.¹⁷² A Q-értékfüggvény konvergenciája Q^* -t jelenti, és kielégíti a Bellman-egyenletet, amelyet a következőképpen határozunk meg:¹⁷³

$$Q^*(s, a) = \mathbb{E}_{s' \sim \varepsilon} [r + \gamma \max_a Q^*(s', a) | s, a]$$

Az π^* optimális π ágens politikája megfelel az Q^* által meghatározott cselekvésnek minden egyes állapotban.¹⁷⁴ Felmerül azonban az a probléma, hogy $Q(s, a)$ értékét minden állapot-akció párra ki kell számítani, ami számítási szempontból kivitelezhetetlen lehet.¹⁷⁵ Például minden olyan állapot-akció pár értékének kiszámítása, ahol a nyers bemenet pixelek egy Atari játékban, óriási számítási teljesítményt igényelne.¹⁷⁶ Az egyik megoldás a Q-érték függvény becslésére egy függvényapproximátor használata:¹⁷⁷

$$Q(s, a; \varnothing) \approx (s, a)$$

Itt \varnothing a függvény paramétereit jelenti.¹⁷⁸ Ha pedig a \varnothing -t egy mély neurális hálózat határozza meg, akkor az algoritmus egy mély megerősítő tanulási algoritmus, amelyet mély Q-hálózatnak ("DQN") nevezünk.¹⁷⁹

A DQN egy mély tanulási modellt, amely egy mély neurális hálózatot ("DNN") kombinál egy Q-tanulási algoritmussal.¹⁸⁰ A DQN a tapasztalatok visszajátszását használja

164 Lásd KOCHENDERFER, *Supra* note at 152,78.

165 Lásd *id.* 79.

166 Sallab et al., 71-72. o. (lásd a 161,71-72. pontot).

167 SUTTON & BARTO, *fenti* 101,108. lábjegyzet.

168 *Id.* 107.

169 862. sz. kérelem, a15, *fenti* megjegyzés 1.

170 Li, Johnson, & Yeung *fenti* jegyzet. 103.

171 SUTTON & BARTO, *fenti* 101,47. lábjegyzet.

172 Sallab et al., *Supra* note at 161,72.

173 862. sz. kérelem, a15, *fenti* megjegyzés 5.

174 MAXIM LAPAN, DEEP REINFORCEMENT LEARNING HANDS-ON, (1022018).

175 Li, Johnson, & Yeung, *Supra* note. 103.

176 LAPAN, *Supra* note at 174,125.

177 862. sz. kérelem, a15, *fenti* megjegyzés 5.

178 *Id.*

179 Li, Johnson, & Yeung, *Supra* note. 103.

180 Manon Legrand, *Deep Reinforcement Learning for Autonomous Vehicle Control Among Human Drivers*, at (2016-17-es 26tanév) (kiadatlan doktori értekezés, Université Libre de Bruxelles) https://ai.vub.ac.be/sites/default/files/thesis_legrand.pdf.

az algoritmus régi tapasztalatainak tárolása a neurális hálózat képzéséhez.¹⁸¹ Egy tapasztalat egy megfigyelt állapot-akció párosból, a kapott azonnali jutalomból és a következő megfigyelt állapotból áll.¹⁸² "Egy ágens tapasztalatát egy t időlépésnél et -vel jelöljük, és egy olyan tuple (s_t, a_t, r_t, s_{t+1}) , amely az aktuális s_t állapotból, az a_t választott akcióból, az r_t jutalomból és a következő s_{t+1} állapotból áll".¹⁸³ Az összes időlépés tapasztalatait egy visszajátszási memóriában tároljuk, sok epizódon keresztül, és a DNN betanításához használjuk.¹⁸⁴ A DNN kimenete egy érvényes akciónak felel meg, mivel a DNN a Q-érték függvény közelítőjeként szolgál.¹⁸⁵ Így a hálózat egy előrecsatolásos áthaladása után a kimenetek az állapot-akció pár becsült Q-értékei.¹⁸⁶ Ez lehetővé teszi, hogy az algoritmus általánosíthatson a múltbeli tapasztalatok összegyűjtött adataiból.¹⁸⁷ Max Tegmark, az MIT professzora szerint ugyanis "a mély megerősítéses tanulás egy teljesen általános technika".¹⁸⁸

III. A MESTERSÉGES INTELLIGENCIA SZABÁLYOZÁSA

A mesterséges intelligencia fenyegetésével foglalkozó jogi tudomány két különálló táborra oszlik.¹⁸⁹ Az egyik tábor elismeri a mesterséges intelligencia rosszindulatú és meggondolatlan használata által jelentett potenciális veszélyeket, míg a másik szerint a mesterséges intelligencia apokalipszise csupán a tudományos fantasztikum szüleménye.¹⁹⁰ E táborok egyike sem foglalkozik igazán a mesterséges intelligencia által jelentett egzisztenciális fenyegetésekkel a katasztrófa megelőzéséhez szükséges azonnalisággal.¹⁹¹ Ez a rész először azokat az érveket tárgyalja, amelyek azzal az elképzeléssel kapcsolatosak, hogy a mesterséges intelligencia nem jelent veszélyt az emberiségre. Ezután ez a rész azokat az érveket tárgyalja, amelyek az AI-szabályozásban az iparági vezetők aggályait tükrözve fejlesztették a tudományt. Végül ez a rész a mesterséges intelligencia szabályozásának három folyamatban lévő és megválaszolatlan kérdésével foglalkozik.

Azok a tudósok, akik azt állítják, hogy a mesterséges intelligencia apokalipszise csupán tudományos fikció, tévednek.¹⁹² Ezeket a tudósokat különösen John McGinnis képviseli, aki megjegyzi, hogy "az emberek által irányíthatatlanná váló gépektől való egzisztenciális rettegés és a politikai aggodalom a gépek pusztító ereje miatt a forradalmasított csatatéren" túlzó.¹⁹³ McGinnis az AGI-vel kapcsolatos problémákat valójában egy gondolkodási hibának tulajdonítja, amikor az emberek antropomorfizálják a mesterséges intelligenciát, és tévesen attól tartanak, hogy az AGI szükségszerűen az emberi rosszindulatot fogja tükrözni.¹⁹⁴ McGinnis tehát a barátságos mesterséges intelligencia lehetőségét sugallja, és arra ösztönöz, hogy hagyjunk fel azzal a feltételezéssel, hogy a mesterséges intelligencia

181 HADO VAN HASSELT, ARTHUR GUEZ, & DAVID SILVER, ASS'N FOR THE ADVANCEMENT OF ARTIFICIAL INTELLIGENCE, PROCEEDINGS OF THE THIRTIETH AAAI CONFERENCE ON ARTIFICIAL INTELLIGENCE: DEEP REINFORCEMENT LEARNING WITH DOUBLE Q-LEARNING (2016), 2094.

182 Volodymyr Mnih et al., *Human-Level Control Through Deep Reinforcement Learning*, NATURE 518 INT'L J. SCI. 529,529 (2015).

183 Legrand, *Supra* note at 180, A 72. "tuple" egy lista vagy tömbhöz hasonló adattárolási formátum. *Id.* 9.

184 VAN HASSELT, GUEZ, & SILVER, *fenti* megjegyzés. 181.

185 Legrand, *Supra* note at 180,27.

186 Mnih et al., *Supra* note. 182.

187 KOCHENDERFER, *supra* note at 152,124.

188 TEGMARK, *Supra* note at 5,85.

189 Lásd Calo, *Supra* note, 6,432. o.; lásd még Matthew U. Scherer, *Regulating Artificial Intelligence Systems*: 29 HARV. J.L. & TECH. 353 (2016).

190 Lásd általában McGinnis, 99. lábjegyzet; lásd még Scherer, 99. lábjegyzet. 189,394.

191 Lásd általában BOSTROM, *Supra* 22. lábjegyzet, at 85.

-
- 192 *Lásd id.*; lásd még TEGMARK, *fenti* megjegyzés. 5.
193 *Lásd* McGinnis, *Supra* note at 99,1254.
194 *Id.*

akaraterővel kell rendelkeznie, mint egy embernek.¹⁹⁵ Feltételezi, hogy az akaraterő hiánya semmissé teszi a gonosz mesterséges intelligenciát övező félelmet.¹⁹⁶

Érdekes módon az antropomorf érv mindkét irányba hat. Nick Bostrom és Eliezer Yudkowsky ugyanis meggyőzően hozta fel az antropomorf érvet annak magyarázatára, hogy az emberek miért fogják drasztikusan alábecsülni a mesterséges intelligencia fejlődését.¹⁹⁷ Bostrom és Yudkowsky azzal érvel, hogy az emberi antropomorfizmus miatt a közvélemény úgy fogja érzékelni, hogy a mesterséges intelligencia fejlődése gyors ütemű lesz.¹⁹⁸ A mesterséges intelligencia emberi antropomorfizmusa itt is arra utal, hogy nem emberi entitásoknak emberi szintű intelligenciát tulajdonítanak.¹⁹⁹ Amint azt az I. részben Einstein és a falu bolondja összehasonlítása is mutatja, az intelligenciaszintek közötti különbség nagyobb relatív léptékben *jelentéktelen*.²⁰⁰ Így az AI fejlődése az AGI és a szuperintelligencia felé a vártnál gyorsabb lesz, mivel a két intelligenciaszint közötti különbség szélesebb léptékben sokkal kisebb, mint azt az emberek felismerik.²⁰¹

Egy másik jogtudós, Ryan Calo még nyersebben érvel amellelt, hogy a mesterséges intelligencia nem jelent egzisztenciális fenyegetést az emberiségre, és hogy az AGI csupán a "képregények anyagát" jelenti.²⁰² Calo továbbá azt állítja, hogy "az AI apokalipsziszre fordított aránytalan figyelem és erőforrások elvonhatják a döntéshozók figyelmét az AI közvetlenebb ártalmainak kezelésétől". . .²⁰³ Azt állítja, hogy a gépi tanulás területén semmi sem utal arra, hogy az emberiség hamarosan képes lesz modellezni az emlősök, nem is beszélve az emberi intelligenciáról.²⁰⁴ Ez az állítás nyilvánvalóan téves. A megerősítő tanulás és a Markov-döntési folyamatok ugyanis szó szerint modellezik az emberi kognitív funkciókat, a döntéshozatalt, a racionális cselekvőképességet és az intelligenciát.²⁰⁵ Ráadásul az adattermelés, a számítási teljesítmény és a globális GDP exponenciális növekedése mind azt a következtetést támasztja alá, hogy az AGI hamarabb érkezik, mint az emberek gondolnák.²⁰⁶

Ezért az AGI által az emberiségre jelentett egzisztenciális fenyegetés azonnali kárt jelent. Ez a fenyegetés nem analóg egy terminátorszerű robot világalomra törésével. Ehelyett ez a fenyegetés a mély megerősítéses tanulóval tanuló ágensekből kifejlesztett AGI terméke.²⁰⁷ Amint az AGI szintű ágensek létrejönnek, gyorsan képesek lesznek arra, hogy szoftverarchitektúrájukat hatékonyabban fejlesszék, mint bármelyik ember. Ezek az ágensek képesek lesznek bármilyen, jutalmazási rendszerrel korrelált cél elérésére egy virtuális környezetben. A mély megerősítő tanulási rendszerek már képesek rakéták, rakéták, autók és repülőgépek irányítására.²⁰⁸ És az ezekhez szükséges szoftverek

195 *Id.* 1263-64.

196 *Id.*

197 BOSTROM, *Supra* note, 22,85. o.; lásd YUDKOWSKY, *Supra* note, 75, 85. o. 21.

198 BOSTROM, *Supra* note at 22,85.

199 YUDKOWSKY, *Supra* note at 75,21.

200 *Id.*

201 BOSTROM, *Supra* note at 22,86.

202 Calo, *Supra* note at 6,432.

203 *Id.* 431.

204 *Id.* 432.

205 Lásd még: TEGMARK, *fentebb*, 5,85. o.; lásd általában: BOSTROM, *fentebb*, 85. o. 22,239.

206 BOSTROM, *Supra* 22. lábjegyzet, 1-4. pont.

207 C-SPAN, *supra* note 4.

208 Mnih et al., *Supra* note, 182,529. o.; lásd még Legrand, *Supra* note, 180,27. o.; lásd még Kendall et al, *supra* note ;149 lásd még: U.S. Patent No. to 8,678,321Bezoz et al. (20145., március).

alkalmazások nyílt forráskódúak.²⁰⁹ Így ma már mindenki, aki rendelkezik internet-hozzáféréssel, potenciálisan hozzáférhet a világ legfejlettebb fegyverzetirányítási rendszereihez is.²¹⁰

A jogi tudomány mégis teljesen figyelmen kívül hagyja ezt a megkerülhetetlen igazságot. Néhány jogtudós azonban tett lépéseket a helyes irányba anélkül, hogy kifejezetten foglalkozott volna az AGI szabályozásának kérdésével. Matthew Scherer például azt állítja, hogy a mesterséges intelligencia szabályozásának kiindulópontja egy olyan jogszabály kellene, hogy legyen, amely meghatározza a mesterséges intelligencia szabályozásának általános elveit.²¹¹ Javasolja a mesterséges intelligencia fejlesztési törvényt ("AIDA"), amely létrehozna egy ügynökséget, amelynek feladata a mesterséges intelligencia rendszerek biztonságának tanúsítása lenne.²¹² Az ügynökségnek szabályokat kellene alkotnia a mesterséges intelligencia meghatározására.²¹³ Az AIDA lényege, hogy a mesterséges intelligencia rendszerek biztonságának értékelésével kapcsolatos érdemi feladatot egy független, szakemberekből álló ügynökségre bízna, így a konkrét mesterséges intelligencia rendszerek biztonságáról szóló döntéseket elszigetelné a választási politika által gyakorolt nyomástól.²¹⁴

Más tanulmányok a különböző szabályozási kereteket tárgyalják, amelyeket a mesterséges intelligenciával kapcsolatos kérdések elemzésére lehet alkalmazni, valamint néhány konkrét példát mutatnak be a mesterséges intelligencia szabályozásának problémáira.²¹⁵ A cikk meggyőzően érvel amellett, hogy a mesterséges intelligenciát, "függetlenül a benne rejlő lehetőségektől, óvatosan kell kezelni".²¹⁶ A szerzők az innováció előmozdítása érdekében árnyalt, érzékeny és alkalmazkodó szabályozási keretet szorgalmaznak.²¹⁷ Miközben az AI-szabályozás területén korlátozott előrelépés történt, az intelligens gépek etikájával és biztonságával kapcsolatos kutatások gyors növekedése nem hozott valódi előrelépést.²¹⁸ Ahogy az egyik írás megjegyzi, "[a] publikált tanulmányok nagy többsége nem tesz mást, mint arról vitatkozik, hogy a meglévő etikai iskolák közül, amelyek évszázadokon át az emberi társadalom igényeinek kielégítésére épültek, melyik lenne a helyes, ha mesterséges utódainkra alkalmaznánk".²¹⁹ Továbbá, még a terület progresszívebb tudományos munkássága is kvázi kizárólag a szűk értelemben vett mesterséges intelligenciára összpontosít, nem pedig az AGI-re.²²⁰ Így a tudósok által javasolt szabályozási keretek egyike sem foglalkozott megfelelően az AGI fejlődésének számos fontos kérdésével.

209 Üdvözljük a *Spinning Up in Deep RL: User Documentation*, OPENAI (utolsó látogatás: 2019. március 20.) spinningup.openai.com.

210 *TensorFlow 2.0 Alpha is Available*, TENSORFLOW, (2019), <https://www.tensorflow.org/install>; lásd még RICHARD WU ET AL., AAAI 2017 FALL SYMPOSIUM SERIES, A FRAMEWORK USING MACHINE VISION AND DEEP REINFORCEMENT LEARNING FOR SELF-LEARNING MOVING OBJECTS IN A VIRTUAL ENVIRONMENT (2017), <https://aaai.org/ocs/index.php/FSS/FSS17/paper/view/16003/15319>.

211 Scherer, *supra* at 189,394.

212 *Id.* 393.

213 *Id.* 394.

214 *Id.* 393.

215 Lásd általában Guihot et al., *Supra* note. 91.

216 *Id.* 454.

217 *Id.*

218 *Lásd* Roman Yampolskiy & Joshua Fox, *Safety Engineering for Artificial General Intelligence*,
TOPOI (322172012).
219 *Id.*
220 *Id.*

A. AZ AGRÁR- ÉS VIDÉKFEJLESZTÉS KORTÁRS KÉRDÉSEI

A szakértők gyanúja szerint a kibertámadók hamarosan olyan stratégiákat fognak alkalmazni, amelyek mély megerősítő tanulási ágenseket használnak olyan támadások kidolgozására, amelyeket a jelenlegi technikai védelmi rendszerek nem képesek megakadályozni.²²¹ Egy tudós konkrétan részletezi a rosszindulatú mesterséges intelligencia szoftverek fejlesztésére vonatkozó irányelveket.²²² Az ösztöndíj azért íródott, hogy bemutassa, hogy gyakorlatilag lehetséges olyan gépi tanulási algoritmusokat fejleszteni, amelyek képesek ártani az embereknek.²²³ Ráadásul az embereknek már ma is megvan a hatalmuk, hogy nukleáris fegyverek használatával elpusztítsák az életet a Föld bolygón, és egy AGI minden bizonnyal ugyanezzel a képességgel rendelkezne.²²⁴ A modern mesterséges intelligenciával foglalkozó tudósok az AGI-k, különösen a mély megerősítő ágensek létrehozásának folyamatát a nukleáris fegyverek építéséhez hasonlítják.²²⁵ Ez a rész három olyan konkrét kérdéssel foglalkozik, amelyeket a mesterséges intelligenciára vonatkozó megfelelő szabályozási keretnek figyelembe kellene vennie.

Az első kérdés a verseny problémája. Ha a szabályozó hatóságok megpróbálják felügyelni az AGI-t fejlesztő vállalatokat, akkor ez a felügyelet elfojtja az innovációt, és lehetővé teszi, hogy olyan országok, mint Kína és Oroszország, az Egyesült Államok előtt fejlesszenek AGI-t.²²⁶ Valójában nagy a valószínűsége annak, hogy bármelyik AGI-t létrehozó szervezet döntő előnyre tesz szert a világ többi részével szemben.²²⁷

A DQN algoritmusokat például gyakran használják részvények kereskedelmére, ahol egy ügynök minden egyes állapotban képes egy részvény megvásárlására, eladására vagy tartására.²²⁸ Az ágens célja a portfólió értékének maximalizálása.²²⁹ A DQN-algoritmusok portfóliókezelésre való alkalmazása sikeresnek bizonyult.²³⁰ Ha egy entitás képes lenne AGI-t létrehozni, akkor ez felhasználható lenne egy olyan ágens létrehozására, amely képes lenne a piacokat úgy manipulálni, hogy egyetlen szereplő minimális idő alatt rendkívüli vagyona tegyen szert.²³¹ Ez lehetővé tenné egy ilyen entitás számára, hogy a tömegek által nem ismert, egységes központi hatalmi erővé fejlődjön.²³²

A verseny problémája még ijesztőbbé válik, ha figyelembe vesszük, hogy a mesterséges intelligencia területén ma a legnagyobb szereplők tőzsdén jegyzett vállalatok. Az olyan vállalatok, mint a Google, a Facebook, az Apple és a Microsoft a mesterséges intelligencia fejlesztésének legnagyobb szereplői közé tartoznak, és technológiájuk teljesítménye és skálázhatósága gyorsan növekszik.²³³ Az e vállalati szereplők, a külföldi kormányok és az Egyesült Államok közötti hatalmi egyenlőtlenség további problémákat vet fel a szabályozók számára.²³⁴ Ha a szövetségi kormány elkezd szabályozni a mesterséges intelligenciát, óvakodnia kell attól, hogy a fejlődés ütemének beföldi lassítása minden bizonnyal hátrányos helyzetbe hozza az Egyesült Államokat a külföldi szereplőkkel szemben. A végső kérdés a

221 BRUNDAGE ET AL., *supra* note 5, at 34; *see also* HYRUM S. ANDERSON, ET AL., LEARNING TO EVADE STATIC PE MACHINE LEARNING MALWARE MODELS VIA REINFORCEMENT LEARNING, (20218) (accessed at <https://arxiv.org/abs/1801.08917>).

222 Federico Pistono & Roman Yampolskiy, *Etikátlan kutatás: (2016)*, (publikálatlan cikk) (elérhető: <https://arxiv.org/abs/1605.02817>).

223 *Id.*

224 Yampolskiy & Fox, *fenti* megjegyzés. 218.

225 BOSTROM, *Supra* note at 22,104.

226 TEGMARK, *Supra* note at 5,9.

227 BOSTROM, *Supra* note at 22,103.

228 LAPAN, *Supra* note at 174,217.

229 *Id.* 220.

230 *Lásd általában* Zhipeng Liang et al., *Deep Reinforcement Learning in Portfolio Management*, (Aug. 201829), (publikálatlan tanulmány) (elérhető a <https://arxiv.org/abs/1808.09940> oldalon).

231 TEGMARK, *fenti* megjegyzés, 5,15-16. o.

232 *Id.*

233 Guihot et al., *supra* note at 91,437.

234 *Id.*

A verseny problémája az, hogy a szabályozóknak a biztonság és a szabadság közötti egyensúlyozással kell szembenéznük. Ha a szabályozók nagyobb hangsúlyt fektetnek a biztonságra, azt a szabadság rovására teszik, amely lehetővé tette a hazai technológiai iparágak vezetőinek az innovációt. Másrészt, ha a szabályozók nagyobb hangsúlyt fektetnek a szabadságra, azt a választópolgárok biztonságának rovására teszik. Ezért a szabályozóknak olyan keretet kell kialakítaniuk, amely érzékeny a vállalatok, a külföldi kormányok és a nemzetbiztonsági ügynökségek közötti versenyre.

A második probléma, amellyel a szabályozó hatóságok szembesülnek, a "magányos farkas" koncepció, amelyben a fenyegetést elszigetelt incidensnek tekintik, nem pedig széles körű társadalmi problémának. A mesterséges intelligencia szabályozása sok szempontból a matematika vagy a számítástechnika szabályozásához hasonló. A mesterséges intelligencia kutatásához ugyanis csak egy személyi számítógépre van szükség.²³⁵ Érdekes módon a tudósok megosztottak az AGI kifejlesztésére irányuló potenciális projekt méretét illetően.²³⁶ Az egyik tudós megjegyzi, hogy az AGI-hez vezető út egy hatalmas kormányzati projekt részeként egy kis csoport vagy akár egyetlen egyén munkájából is megvalósulhat.²³⁷ Az AGI-hez vezető út nagyságrendje nagymértékben függ a mesterséges intelligencia eléréséhez használt módszerektől.²³⁸ Például, ha a teljes agy emulációjának jelenlegi módszereit alkalmazzák, akkor valószínű, hogy az AGI kifejlesztéséhez hatalmas mennyiségű kódot kell kidolgozniuk a szakértő informatikusoknak és mérnököknek.²³⁹ Fontos megjegyezni, hogy míg maga a projekt hatalmas léptékű lehet, addig a mesterséges intelligenciától az AGI-ig való áttörés megvalósításával megbízott egyéni csoport nagyon kicsi lehet.²⁴⁰ A Manhattan-projekt például csúcspontján nagyjából 130 000 embert foglalkoztatott.²⁴¹ Az atombombát azonban tudósok és mérnökök egy kisebb csoportja hozta létre, J. Robert Oppenheimer és Leslie Groves tábornok vezetésével a Los Alamos-i Tudományos Laboratóriumban.²⁴²

A magányos farkas problémájának egy másik kérdése akkor fog jelentkezni, ha a mesterséges intelligencia területén egyetlen ember fog áttörést elérni. Ebben az esetben lehetséges, hogy minden, amit jelenleg a mesterséges intelligenciáról tudunk, a süllyesztőbe kerül. A tudománytól nem idegenek az egyszerű, de forradalmi áttörések, amelyek gyökeresen megváltoztatják az emberiség természeti világról alkotott képét.²⁴³ Egy tudós szerint azonban valószínű, hogy a szabályozó szervek a legtöbb, AGI kifejlesztésére potenciálisan képes emberről tudnának.²⁴⁴ Bár meg kell jegyezni, hogy egy olyan megvilágosodás az AI terén, mint amilyen Einstein *Annus Mirabilis*-dokumentumaiban kifejezett fizikai megvilágosodása volt, nem zárható ki a lehetőségek köréből. Ezért lehetséges, hogy egyetlen egyén lesz az első, aki létrehozza az AGI-t, és röviddel ezután példátlan mértékű hatalomra tehet szert.²⁴⁵ A szabályozóknak olyan keretrendszer kell kialakítaniuk, amely lehetővé teszi a magányos AGI támadások és fenyegetések azonosítására és megelőzésére szolgáló technológia alkalmazását.

A harmadik kérdés, amellyel a technológiai szabályozóknak szembe kell nézniük, az ellenőrzési probléma. Az ellenőrzési probléma két különböző módon elemezhető a megbízóügynöki keretrendszerben.²⁴⁶

235 BOSTROM, *Supra* note at 22,103.

236 *Id.*

237 *Id.* 101.

238 *Id.*

239 KURZWEIL, *Supra* note at 55,124.

240 BOSTROM, *Supra* note at 22,101.

241 F.G. GOSLING, A MANHATTAN PROJEKT: AZ ATOMBOMBA KÉSZÍTÉSE (541999).

242 *Id.* 35.

243 Lásd általában Albert Einstein, *On the Electrodynamics of Moving Bodies* (1905) in THE COLLECTED PAPERS OF ALBERT EINSTEIN: THE SWISS YEARS: WRITINGS, 1900-1909 (140Anna Beck trans., 1990),

<https://einsteinpapers.press.princeton.edu/vol2-trans/154>; *lásd még* BRIAN GREENE, FABRIC OF THE COSMOS (1282005).

244 BOSTROM, *Supra* note at 22,103.

245 *Id.*

246 *Id.* 155.

Az első keretrendszer akkor létezik, ha a projekt szponzora megbízóként, a tudósok és mérnökök egy csoportja pedig a projekt szponzorának megbízottjaként jár el.²⁴⁷ Ebben a keretben az ellenőrzési probléma akkor jelentkezik, ha az AGI-t fejlesztő tudósok és mérnökök rosszindulatú célokra használják fel a munkájuk során szerzett ismereteket és információkat.²⁴⁸ Például az Apple, a Google és a Facebook kutatói saját cégeik mesterséges intelligenciafejlesztésének eredményeképpen hatalmas hatalomra²⁴⁹ és arra a képességre tettek szert, hogy fejlett mesterséges intelligencia-rendszereket fejlesszenek ki vagy módosítsanak saját személyes hasznukra vagy mások kárára.

A második keretrendszerben a megbízó az emberi alkotó, az ügynök pedig a mesterséges intelligencia rendszer.²⁵⁰ Ebben a keretrendszerben az irányítási probléma akkor jelentkezik, ha egy AGI-rendszert fejlesszenek ki, és annak cselekvései az alkotója által nem ellenőrizhetők.²⁵¹ Az AGI megfékezésére több különböző módszert mutattak be. Nick Bostrom például bokszozási módszereket javasolt az AGI információhoz való hozzáféréseinek felosztására és korlátozására.²⁵² Max Tegmark továbbá felvetette, hogy lehetséges egy "kapuőr mesterséges intelligencia" létrehozása, egy olyan szuperintelligencia, amelynek célja, hogy a lehető legkisebb mértékben avatkozzon be, hogy megakadályozza egy másik szuperintelligencia létrejöttét.²⁵³ Ezért a szabályozóknak olyan keretet kell kialakítaniuk, amely ellenőrzi, hogy az AI-kutatók hogyan használják fel hatalmukat, és olyan keretet, amely lehetővé teszi az AGI-rendszerek szabályozását, hogy azokat emberi szereplők irányíthassák.

Összefoglalva, a mesterséges intelligenciát szabályozó hatóságok három fő problémája a verseny, a magányos farkas koncepció és az ellenőrzés. E problémák gyakorlati megoldásának egyik gyakorlati módja az önszabályozó AGI-technológia.²⁵⁴ Ehhez lényegében az AGI értékeinek az AGI létrehozójának értékeivel összhangban lévő programozására lesz szükség.²⁵⁵ Ennek a megoldásnak az egyik fő előnye, hogy lehetővé teszi, hogy az állami szabályozók viszonylag sötétben maradjanak az AI-technológia működésével kapcsolatban.²⁵⁶ Ezzel a megoldással kapcsolatban azonban két jelentős probléma van. Először is, ha létezik egy olyan AGI-rendszer, amely képes szabályozni az összes többi AI-rendszert, akkor szükség lesz egy olyan szabályozási mechanizmusra, amely megfékezi a szabályozó AGI-t, hogy az ne váljon az emberek felett ellenőrzést gyakorló egységes hatalommá. Másodsor, az emberiségnek biztosítania kell, hogy a szabályozó AGI-t ne győzhesse le és ne győzhesse le semmilyen más mesterséges intelligencia vagy AGI. Valóban, ha létezne egy AGI, amely képes fejleszteni önmagát, akkor bármely emberi programozó képességei gyorsan messze lemaradnának, miközben Irving J. Good hírhedt szavai: "...az utolsó találmány, amit az embernek valaha is meg kell tennie..." profétai jelleggel fognak visszhangozni.²⁵⁷

KÖVETKEZTETÉS

A tudományban a lehetőségek spektruma van lefektetve. Az egyik oldalon azok állnak, akik szerint a mesterséges intelligencia a közeljövőben örökre meg fogja változtatni az emberi életet, a másik oldalon pedig azok, akik szerint a mesterséges intelligencia apokalipszise csupán tudományos fikció. Az igazság az, hogy

247 *Id.* 155-56.

248 *Id.*

249 Guihot et al., *supra* note at 91,455.

250 BOSTROM, *Supra* note at 22,156.

251 TEGMARK, *Supra* note at 5,187.

252 BOSTROM, *Supra* note at 22,156.

253 TEGMARK, *Supra* note at 5,176.

254 Guihot et al., *Supra* note, 91,433-37. o.

Elektronikusan elérhető a következő címen:

<https://ssrn.com/abstract=3261254>

-
- 255 TEGMARK, *Supra* note at 5,261.
256 BOSTROM, *Supra* note at 22,156.
257 Good, *fenti megjegyzés* 65.

hogy egyik tábor sem érti teljesen az AGI-t vagy annak a világunkra gyakorolt hatását.²⁵⁸ 1988-ban bekövetkezett halálakor Richard Feynman Nobel-díjas fizikus tábláján a következő szavak szerepeltek: "Amit nem tudok létrehozni, azt nem értem".²⁵⁹ Ebből következik, hogy amíg az emberiség nem hozza létre az AGI-t, addig az emberiség számára felfoghatatlan. Ez a valóság ironikus sorsot jelent az emberiség számára. Az emberiségnek ugyanis először meg kell értenie az AGI-t ahhoz, hogy irányítani tudja, de az emberiség nem értheti meg az AGI-t, amíg nem hozza létre. Továbbá Max Tegmark szerint "fogalmunk sincs, mi fog történni, ha az emberiségnek sikerül emberi szintű AGI-t létrehoznia".²⁶⁰ Így nem vehetjük biztosra, hogy az AGI létrehozása pozitív kimenetelű lesz.²⁶¹ A mesterséges intelligencia területén ugyanis általános konszenzus van abban, hogy egyetlen szabályrendszer sem képes kontrollálni mindazt, amit az AGI szabályozása megkövetel.²⁶² De bármennyi mintát is lehet felismerni és trendeket követni, a mesterséges intelligencia jövője nem fog magától megtörténni.²⁶³

A legtöbb ember azt gondolja, hogy a múlt determinisztikus kapcsolatban áll a jövővel, de az igazság az, hogy a jövő alapvetően bizonytalan.²⁶⁴ A huszadik század eleje óta az emberiség meggyőző bizonyítékokkal rendelkezik arról, hogy az emberek által a mindennapi tapasztalatainkban érzékelt téridő egésze csak az egyén szubjektív megfigyeléséhez képest létezik.²⁶⁵ A kvantumfizikában, valamint a fekete lyukak középpontjában pedig a klasszikus fizika és a téridő törvényei teljesen összeomlanak.²⁶⁶ Ez azért fontos, mert a létezés alapvető szintjén a téridő nélkül a masszív részecskék függetlensége elpárolog, és az ember által érzékelt idő előrehaladása megszűnik létezni.²⁶⁷ Ezt bizonyítják a szuperpozíció, a nem-lokalitás, az időszimmetria és a kvantum összefonódás elvei.²⁶⁸ Valójában a jövő, csakúgy, mint a múlt, a kvantumbizonytalanság alapvető állapotában létezik. A mesterséges intelligencia a kulcsa annak, hogy az ilyen bizonytalanság mellett is maximalizáljuk a jólét valószínűségét. Ezért azonnal dolgoznunk kell a biztonságos és virágzó jövő megteremtésén.²⁶⁹

258 Lásd BOSTROM, *fenti* megjegyzés (azzal 22érvvel, hogy már egy szimulációban élünk; ez egy valós lehetőség, és erős érv, hogy egy entitás, esetleg a téridőben, megérti az AGI-t).

259 Lásd Michael Way, "Amit nem tudok létrehozni, azt nem értem", J. OF CELL SCI. (2017), <http://jcs.biologists.org/content/joces/130/18/2941.full.pdf>.

260 Lásd TEGMARK, *fenti* jegyzet, lásd:5, TEGMARK, *fenti* jegyzet. 156.

261 Lásd PETER THEIL, ZERO TO ONE (1952014).

262 Lásd Yampolskiy & Fox, *Supra* note. 218.

263 Lásd THEIL, *fenti* megjegyzés 261.

264 Lásd általában Einstein, *Supra* note 243.

265 Lásd általában *id.*

266 Lásd STEPHEN HAWKING, A BRIEF HISTORY OF TIME 111 (1996); lásd még ROBERT J. SPITZER, NEW PROOFS FOR THE EXISTENCE OF GOD (1232010).

267 Lásd GREENE, *fenti* jegyzet, 243,192-93.

o. 268 Lásd *id.* 199-208.

269 Lásd THEIL, *fenti* megjegyzés 261.

A függelék

A jelölés összefoglalása	
Jelölés	Jelentése ²⁷⁰
$Q^*(s, a)$	Az a cselekvés értéke az optimális politika szerint
γ	Kedvezménytényező
$E[x]$	Véletlen változó várakozása
$\arg \max_a f(a)$	Az a olyan értéke, amelynél $f(a)$ maximális értéket vesz fel.
r	Jutalom
s_t	Állapot t időpontban
π	Politika
π^*	Optimális politika
$V\pi(s)$	A politika végrehajtásának várható értéke egy adott államból.

270 Lásd általában SUTTON & BARTO, *fenti* megjegyzés. 101.