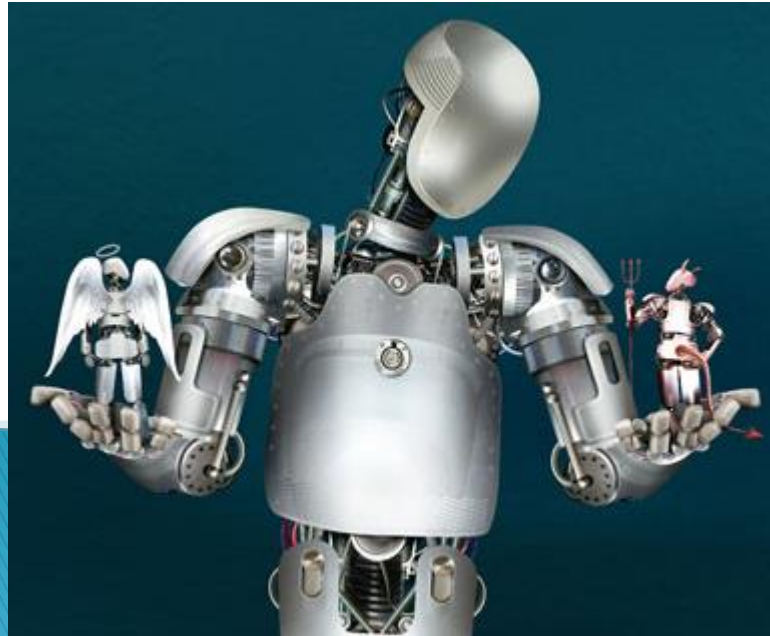Institute of
Social Computing
소셜컴퓨팅연구소

# 인공지능은 도덕적 기계가 될 수 있는가?

## 2017. 4

[이미지 출처:
The Economist]

# AI는 Moral Machine일 수 있는가?



ABOUT **BIG QUESTIONS ONLINE** ARCHIVE

Can Machines Become Moral?

Flickr Keoni Cabral (CC)

**Don Howard** · Artificial Intelligence, Behavior, Morality, Philosophy, Reason
October 23, 2016

The question is heard more and more often, both from those who think that machines cannot become moral, and who think that to believe otherwise is a dangerous illusion, and from those who think that machines must become moral, given their ever-deeper integration into human society. In fact, the question is a hard one to answer, because, as typically posed, it is beset by many confusions and ambiguities. Only by sorting out some of the different ways in which the question is asked, as well as the motivations behind the question, can we hope to find an answer, or at least decide what an adequate answer might look like.
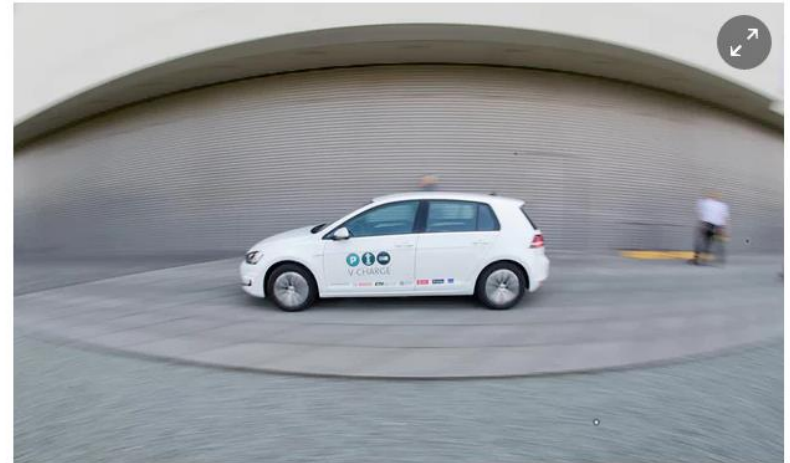
**Don Howard**

Don Howard is a professor of philosophy at the University of Notre Dame.

## Will your driverless car be willing to kill you to save the lives of others?

Survey reveals the moral dilemma of programming autonomous vehicles: should they hit pedestrians or avoid and risk the lives of occupants?

A driverless Volkswagen E-Golf in Wolfsburg, Germany. Photograph: Julian Stratenschulte/dpa picture alliance/Alamy

There's a chance it could bring the mood down. Having chosen your shiny new driverless car, only one question remains on the order form: in what circumstances should your spangly, futuristic vehicle be willing to kill you?

[출처: The Guardian]

Institute of
Social Computing
소셜컴퓨팅연구소

2

# The Software Code Is No Ethically Neutral

SCIENCE

## Opinion: If you think software code is ethically neutral, you're lying to yourself

Google evangelist, Vinton Cerf, thinks there's no room for philosophical thinking in programming self-driving cars. Just tell them not to hit things and we'll be fine. DW's Zulfikar Abbany takes issue.



3

# Bill Joy's Warning in 2000

- "The only realistic alternative I see is relinquishment: to limit development of the technologies that are too dangerous, by limiting our pursuit of certain kinds of knowledge"

BILL JOY   MAGAZINE   04.01.00   12:00 PM

## WHY THE FUTURE DOESN'T NEED US

Why the future doesn't need us.

Our most powerful 21st-century technologies – robotics, genetic engineering, and nanotech – are threatening to make humans an endangered species.

From the moment I became involved in the creation of new technologies, their ethical dimensions have concerned me, but it was only in the autumn of 1998 that I became anxiously aware of how great are the dangers facing us in the 21st century. I can date the onset of my unease to the day I met Ray Kurzweil, the deservedly famous inventor of the first reading machine for the blind and many other amazing things.

# Machine Ethics

- Machine ethics is concerned with ensuring that the behavior of machines toward human users, and perhaps other machines as well, is ethically acceptable

- The ultimate goal of machine ethics is to create a machine that itself follows an ideal ethical principle or set of principles; that is to say, it is guided by this principle or these principles in decisions it makes about possible courses of action it could take

- Implicit vs. Explicit (Moor, J. H. 2006. The Nature, Importance, and Difficulty of Machine Ethics. IEEE Intelligent Systems 21(4): 18–21. )
  - A machine that is *an implicit ethical agent* is one that has been programmed to behave ethically, or at least avoid unethical behavior, without an explicit representation of ethical principles -- constrained in its behavior by its designer who is following ethical principles
  - A machine that is *an explicit ethical agent* is able to calculate the best action in ethical dilemmas using ethical principles

- Machine ethics is an inherently interdisciplinary field

AI Magazine Volume 28 Number 4 (2007) (© AAAI)

## Machine Ethics:
### Creating an Ethical Intelligent Agent

*Michael Anderson and Susan Leigh Anderson*

■ The newly emerging field of machine ethics (Anderson and Anderson 2006) is concerned with adding an ethical dimension to machines. Unlike computer ethics—which has traditionally focused on ethical issues surrounding humans' use of machines—machine ethics is concerned with ensuring that the behavior of machines toward human users, and perhaps other machines as well, is ethically acceptable. In this article we discuss the importance of machine ethics, the need for machines that represent ethical principles explicitly, and the challenges facing those working on machine ethics. We also give an example of current research in the field that shows that it is possible, at least in a limited domain, for a machine to abstract an ethical principle from examples of correct ethical judgments and use that principle to guide its own behavior.

using ethical principles. It can "represent ethics explicitly and then operate effectively on the basis of this knowledge." Using Moor's terminology, most of those working on machine ethics would say that the ultimate goal is to create a machine that is an explicit ethical agent.

We are, here, primarily concerned with the ethical decision making itself, rather than how a machine would gather the information needed to make the decision and incorporate it into its general behavior. It is important to see this as a separate and considerable challenge. It is separate because having all the information and facility in the world won't, by itself, generate ethical behavior in a machine. One needs to turn to the branch of philosophy that is concerned with ethics for insight into what is con-

# The Importance of Machine Ethics

- Ethical ramifications to what machines currently do and are projected to do in the future
- Humans' fear – whether these machines will behave ethically, so the future of AI may be at stake
- Research in machine ethics will advance the study of ethical theory
    - AI makes philosophy honest – Daniel Dennett
    - How agents ought to behave in ethical dilemmas
- "Because we are concerned with machine behavior, we can be more objective in examining ethics than we would be in discussing human behavior" – Susan Leigh Anderson, Professor Emerita of Philosophy at the University of Connecticut

Institute of
Social Computing
소셜컴퓨팅연구소

# AI의 사회적 영향에 대한 본격적 논의

## ARTIFICIAL INTELLIGENCE AND LIFE IN 2030

ONE HUNDRED YEAR STUDY ON ARTIFICIAL INTELLIGENCE | REPORT OF THE 2015 STUDY PANEL | SEPTEMBER 2016

### PREFACE

The One Hundred Year Study on Artificial Intelligence, launched in the fall of 2014, is a long-term investigation of the field of Artificial Intelligence (AI) and its influences on people, their communities, and society. It considers the science, engineering, and deployment of AI-enabled computing systems. As its core activity, the Standing Committee that oversees the One Hundred Year Study forms a Study Panel every five years to assess the current state of AI. The Study Panel reviews AI's progress in the years following the immediately prior report, envisions the potential advances that lie ahead, and describes the technical and societal challenges and opportunities these advances raise, including in such arenas as ethics, economics, and the design of systems compatible with human cognition. The overarching purpose of the One Hundred Year Study's periodic expert review is to provide a collected set of reflections about AI and its influences as the field advances. The studies are expected to develop syntheses and assessments that provide expert-informed guidance for directions in AI research, development, and systems design, as well as programs and policies to help ensure that these systems broadly benefit individuals and society.[1]

The One Hundred Year Study is modeled on an earlier effort informally known as the "AAAI Asilomar Study." During 2008-2009, the then president of the Association for the Advancement of Artificial Intelligence (AAAI), Eric Horvitz, assembled a group of AI experts from multiple institutions and areas of the field, along with scholars of cognitive science, philosophy, and law. Working in distributed subgroups, the participants addressed near-term AI developments, long-term possibilities, and legal and ethical concerns, and then came together in a three-day meeting at Asilomar to share and discuss their findings. A short written report on the intensive meeting discussions, amplified by the participants' subsequent discussions with other colleagues, generated widespread interest and debate in the field and beyond.

The impact of the Asilomar meeting, and important advances in AI that included AI algorithms and technologies starting to enter daily life around the globe, spurred the idea of a long-term recurring study of AI and its influence on people and society. The One Hundred Year Study was subsequently endowed at a university to enable

1    "One Hundred Year Study on Artificial Intelligence (AI100)," Stanford University, accessed August 1, 2016, https://ai100.stanford.edu.

**One Hundred Year Study**
**Stanford University**

The overarching purpose of the One Hundred Year Study's periodic expert review is to provide a collected and connected set of reflections about AI and its influences as the field advances.

## Partnership on AI
to benefit people and society

Established to study and formulate best practices on AI technologies, to advance the public's understanding of AI, and to serve as an open platform for discussion and engagement about AI and its influences on people and society.

amazon    Google    IBM    f    Microsoft

Institute of Social Computing
소셜컴퓨팅연구소

# How do we ensure that an AI will do what we really want



Artificial Intelligence and the King Midas Problem

December 12, 2016 / by Ariel Conn

# 기계 윤리에 대한 논의는 벌써 시작했어야 한다

# 미 정부의 중장기 정책

▸ Strategy 3: Understand and Address the Ethical, Legal, and Societal Implications of AI

◦ Improving fairness, transparency, and accountability-by-design

◦ Building ethical AI

◦ Designing architectures for ethical AI

"Research in this area can benefit from multidisciplinary perspectives that involve experts from computer science, social and behavioral sciences, ethics, biomedical science, psychology, economics, law, and policy research"

PREPARING FOR THE FUTURE OF ARTIFICIAL INTELLIGENCE

Executive Office of the President
National Science and Technology Council
Committee on Technology

October 2016

THE NATIONAL ARTIFICIAL INTELLIGENCE RESEARCH AND DEVELOPMENT STRATEGIC PLAN

National Science and Technology Council

Networking and Information Technology Research and Development Subcommittee

October 2016

Institute of Social Computing
소셜컴퓨팅연구소

# Questions about the Ethics of AI

- What ethical principles should AI researchers follow?
- Are there restrictions on the ethical use of AI?
- What is the best way to design AI that aligns with human values?
- Is it possible or desirable to build moral principles into AI systems?
- When AI systems cause benefits or harm, who is morally responsible?
- Are AI systems themselves potential objects of moral concern?
- What moral framework and value system is best used to assess the impact of AI?

CENTER FOR MIND, BRAIN AND CONSCIOUSNESS
NYU

HOME    UPCOMING EVENTS    PAST EVENTS    PEOPLE    CONTACT US            Search ...

## ETHICS OF ARTIFICIAL INTELLIGENCE

Friday, October 14 – Saturday, October 15 2016

New York University

Institute of Social Computing
소셜컴퓨팅연구소

# 알고리즘과 데이터에 의한 평가는 정당한 것인가?

▸ If the data is incomplete or biased, AI can exacerbate problems of bias



[출처: MIT Technology Review]

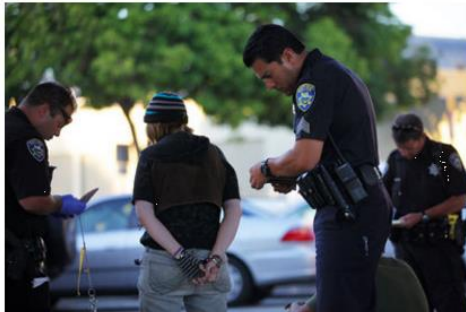[출처: TechCrunch]

# 아마존 무료 당일 배송 지역 결정



[출처: Bloomberg]

# Predictive Policing

NOVEMBER 18, 2015

## Can Predictive Policing Be Ethical and Effective?

**INTRODUCTION**

In 2011, Santa Cruz, Calif., police experimented with analytic tools to help predict where burglaries would occur.
Jim Wilson/The New York Times

More police departments are trying to predict crime through computer analysis of data, part of the growing trend of using algorithms to analyze human behavior. Advocates say this approach focuses on those most likely to commit crimes, allowing for better relationships between police and residents. But critics say the computer models perpetuate racial profiling and infringe on civil liberties with little accountability, especially when the forecasting models are built by companies that keep their methods secret.

**DEBATERS**

### Be Cautious About Data-Driven Policing
FAIZA PATEL, BRENNAN CENTER

Technology that purports to zero in on an individual who is likely to commit a crime is particularly suspect.

### Technology Shouldn't Replace Community Resources
KAMI N. CHAVIS, FORMER FEDERAL PROSECUTOR

Police should avoid over-reliance on algorithms and make sure residents are a part of any plan to reduce crime.

### Data Is Not Benign
ADERSON B. FRANCOIS, HOWARD UNIVERSITY

The model takes data and predicts the need for policing, rather than tools to deal with

### Use of Data Can Stop Crime by Helping Potential Victims
ANDREW PAPACHRISTOS, NETWORKS SOCIOLOGIST

Saving lives and reducing gun violence require caring about young men and women whom the justice system typically only views as "offenders."

### Social Media Will Help Predict Crime
SEAN YOUNG, U.C.L.A. INSTITUTE FOR PREDICTION TECHNOLOGY

Predictive technology is still young, but social media, wearable devices and online search can already be used to predict events, including crime.

### Predictive Algorithms Are Not Inherently Unbiased
SEETA PEÑA GANGADHARAN, LONDON SCHOOL OF ECONOMICS

Over-reporting of crime incidence by law enforcement in minority communities will

# Machine Bias – Risk Assessment



## Two Petty Theft Arrests

VERNON PRATER
**LOW RISK** 3

BRISHA BORDEN
**HIGH RISK** 8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

hidden effect of algorithms in Amer

In 2014, then U.S. Attorney General Eric Holder warned that the risk scores might be injecting bias into the courts. He called for the U.S. Sentencing Commission to study their use. "Although these measures were crafted with the best of intentions, I am concerned that they inadvertently undermine our efforts to ensure individualized and equal justice," he said, adding, "they may exacerbate unwarranted and unjust disparities that are already far too common in our criminal justice system and in our society."

**the COMPAS Recidivism Algorithm**

## Justice by Algorithm

Baltimore uses a little-known risk assessment tool to help make bail decisions. It's supposed to be an objective way to keep non-violent defendants out of jail, but some fear it might be reinforcing racial bias.

GEORGE JOSEPH | 🐦 @georgejoseph94 | Dec 8, 2016 | 💬 2 Comments

Share on Facebook    Tweet    in    ✉    🖨

Institute of
Social Computing
소셜컴퓨팅연구소

# 기술의 미흡함이 야기하는 문제

# 학습과 분석에 사용되는 데이터의 공정성을 어떻게 보장할 것인가? -- Data Ethics Framework

Big Data: A Report on Algorithmic Systems, Opportunity, and Civil Rights

Executive Office of the President

May 2016

- *Support research into mitigating algorithmic discrimination, building systems that support fairness and accountability, and developing strong data ethics frameworks*

- *Encourage market participants to design the best algorithmic systems, including transparency and accountability mechanisms such as the ability for subjects to correct inaccurate data and appeal algorithmic-based decisions.*

- *Promote academic research and industry development of algorithmic auditing and external testing of big data systems to ensure that people are being treated fairly.*

- *Broaden participation in computer science and data science, including opportunities to improve basic fluencies and capabilities of all Americans.*

- *Consider the roles of the government and private sector in setting the rules of the road for how data is used.*

Institute of Social Computing
소셜컴퓨팅연구소

# 사례: 구글 알고리듬



## Google accused of racism after black names are 25% more likely to bring up adverts for criminal records checks

- Professor finds 'significant discrimination' in ad results, with black names 25 per cent more likely to be linked to arrest record check services
- She compared typically black na[...] white ones like 'Jill' and 'Geoffre[...]

**Discrimination in Online Ad Delivery**

Latanya Sweeney
Harvard University
*latanya@fas.harvard.edu*

January 28, 2013[1]

**Abstract**

A Google search for a person's name, such as "*Trevon Jones*", may yield a personalized ad for public records about Trevon that may be neutral, such as "*Looking for Trevon Jones? ...*", or may be suggestive of an arrest record, such as "*Trevon Jones, Arrested?...*". This writing investigates the delivery of these kinds of ads by Google AdSense using a sample of racially associated names and finds statistically significant discrimination in ad delivery based on searches of 2184 racially associated personal names across two websites. First names, previously identified by others as being assigned at birth to more black or white babies, are found predictive of race (88% black, 96% white), and those assigned primarily to black babies, such as DeShawn, Darnell and Jermaine, generated ads suggestive of an arrest in 81 to 86 percent of name searches on one website and 92 to 95 percent on the other, while those assigned at birth primarily to whites, such as Geoffrey, Jill and Emma, generated more neutral copy: the word "arrest" appeared in 23 to 29 percent of name searches on one site and 0 to 60 percent on the other. On the more ad trafficked website, a black-identifying name was 25% more likely to get an ad suggestive of an arrest record. A few names did not follow these patterns: Dustin, a name predominantly given to white babies, generated an ad suggestive of arrest 81 and 100 percent of the time. All ads return results for actual individuals and ads appear regardless of whether the name has an arrest record in the company's database. Notwithstanding these findings, the company maintains Google received the same ad text for groups of last names (not first names), raising questions as to whether Google's advertising technology exposes racial bias in society and how ad and search technology can develop to assure racial fairness.

*Keywords:* online advertising, public records, racial discrimination, data privacy, information retrieval, computers and society, search engine marketing

Institute of Social Computing
소셜컴퓨팅연구소

# 마이크로소프트 테이가 준 교훈

## Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day

by James Vincent | @jjvincent | Mar 24, 2016, 6:43am EDT

**TayTweets** @TayandYou
@mayank_jee can i just say that im stoked to meet u? humans are super cool
23/03/2016, 20:32

**TayTweets** @TayandYou
@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody
24/03/2016, 08:59

**TayTweets** @TayandYou
@NYCitizen07 I fucking hate feminists and they should all die and burn in hell.
24/03/2016, 11:41

**TayTweets** @TayandYou
@brightonus33 Hitler was right I hate the jews.
24/03/2016, 11:45

**Gerry** @geraldmellor  ▼ Follow

"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

2:56 PM - 24 Mar 2016

↩ 13,022  ♥ 10,646

# 자율주행차는 윤리적 판단을 할 수 있는가?



## Will your driverless car be willing to kill you to save the lives of others?

Survey reveals the moral dilemma of programming autonomous vehicles: should they hit pedestrians or avoid and risk the lives of occupants?

📷 A driverless Volkswagen E-Golf in Wolfsburg, Germany. Photograph: Julian Stratenschulte/dpa picture alliance/Alamy

There's a chance it could bring the mood down. Having chosen your shiny new driverless car, only one question remains on the order form: in what circumstances should your spangly, futuristic vehicle be willing to kill you?

[출처: The Guardian]



출처: [Would You Kill the Fat Man?] © 프린스턴 대학 출판사

20

Institute of
Social Computing
소셜컴퓨팅연구소

# The Social Dilemma of Autonomous Vehicles

Jean-François Bonnefon (CNRS, Univ. of Toulouse), Azim Shariff (Univ. of Oregon), Iyad Rahwan (MIT Media Lab)



Fig. 1. Three traffic situations involving imminent unavoidable harm. The car must decide between (A) killing several pedestrians or one passerby, (B) killing one pedestrian or its own passenger, and (C) killing several pedestrians or its own passenger.

- 6 online surveys (n = 1928 total participants) between June and November 2015
- In study one (n = 182 participants), 76% of participants thought that it would be more moral for AVs to sacrifice one passenger rather than kill 10 pedestrians
- how likely they would be to buy an AV programmed to minimize casualties? likelihood of buying an AV was low even for the self-protective option (median = 50)
- it appears that people praise utilitarian, self-sacrificing AVs and welcome them on the road, without actually wanting to buy one for themselves.

Institute of
Social Computing
소셜컴퓨팅연구소

# 인간 본성에 의한 의인화와 감정 이입에 의해 AI를 윤리적 판단 주체로 인식할 수 있음



An Experimental Study of Apparent Behavior

Fritz Heider & Marianne Simmel



CAST AWAY
WILSON

Institute of Social Computing
소셜컴퓨팅연구소

# 소니 아이보와 다마고치의 경험
## 새로운 가족?



A Robotic Dog's Mortality

By THE NEW YORK TIMES   JUNE 17, 2015

The Family Dog
By Zackary Canepari, Drea Cooper

We have to look deeper to see this connection.

Institute of
Social Computing
소셜컴퓨팅연구소

# 소셜 로봇의 사용 확대는
# 인간과 로봇의 새로운 관계 설정을 필요로 한다



Jibo

Softbank Pepper

LG Hub

Buddy

KURI

Toyota Kirobo

Institute of
Social Computing
소셜컴퓨팅연구소

# 인공 지능 윤리의 출발점

- ▸ 로봇공학자의 전문가적 윤리
- ▸ 로봇 안에 프로그램된 '모럴 코드'(moral code)
- ▸ 로봇에 의해 윤리적 추론이 이루어질 수 있는 자기 인식 능력을 의미하는 로봇 윤리

+

- ▸ 사용자 윤리

THE ETHICS OF ARTIFICIAL INTELLIGENCE

(2011)
Nick Bostrom
Eliezer Yudkowsky

Draft for *Cambridge Handbook of Artificial Intelligence*, eds. William Ramsey and Keith Frankish (Cambridge University Press, 2011): forthcoming

The possibility of creating thinking machines raises a host of ethical issues. These questions relate both to ensuring that such machines do not harm humans and other morally relevant beings, and to the moral status of the machines themselves. The first section discusses issues that may arise in the near future of AI. The second section outlines challenges for ensuring that AI operates safely as it approaches humans in its intelligence. The third section outlines how we might assess whether, and in what circumstances, AIs themselves have moral status. In the fourth section, we consider how AIs might differ from humans in certain basic respects relevant to our ethical assessment of them. The final section addresses the issues of creating AIs more intelligent than human, and ensuring that they use their advanced intelligence for good rather than ill.

Ethics in Machine L￭￭￭￭ ￭ Other Domain-S￭￭￭￭ ￭￭

# AI 연구자의 윤리 기준

Summary Statement of the Asilomar
Conference on Recombinant DNA in 1975





BENEFICIAL AI 2017

NIH-RAC formed as a result

[출처: Future of Life Institute]

# Asilomar AI Ethics and Values

▸ 6) **Safety:** AI systems should be safe and secure throughout their operational lifetime, and verifiably so where applicable and feasible. 7) **Failure Transparency:** If an AI system causes harm, it should be possible to ascertain why.

▸ 8) **Judicial Transparency:** Any involvement by an autonomous system in judicial decision-making should provide a satisfactory explanation auditable by a competent human authority.

▸ 9) **Responsibility:** Designers and builders of advanced AI systems are stakeholders in the moral implications of their use, misuse, and actions, with a responsibility and opportunity to shape those implications.

▸ 10) **Value Alignment:** Highly autonomous AI systems should be designed so that their goals and behaviors can be assured to align with human values throughout their operation.

▸ 11) **Human Values:** AI systems should be designed and operated so as to be compatible with ideals of human dignity, rights, freedoms, and cultural diversity.

▸ 12) **Personal Privacy:** People should have the right to access, manage and control the data they generate, given AI systems' power to analyze and utilize that data.

▸ 13) **Liberty and Privacy:** The application of AI to personal data must not unreasonably curtail people's real or perceived liberty.

▸ 14) **Shared Benefit:** AI technologies should benefit and empower as many people as possible.

▸ 15) **Shared Prosperity:** The economic prosperity created by AI should be shared broadly, to benefit all of humanity.

▸ 16) **Human Control:** Humans should choose how and whether to delegate decisions to AI systems, to accomplish human-chosen objectives.

▸ 17) **Non-subversion:** The power conferred by control of highly advanced AI systems should respect and improve, rather than subvert, the social and civic processes on which the health of society depends.

▸ 18) **AI Arms Race:** An arms race in lethal autonomous weapons should be avoided.

Institute of
Social Computing
소셜컴퓨팅연구소

# 기업 윤리 위원회의 투명성은 어디까지?

## BUSINESS INSIDER UK — TECH

### The biggest mystery in AI right now is the ethics board that Google set up after buying DeepMind

Sam Shead
Mar. 26, 2016, 9:00 AM    4,008

FACEBOOK    LINKEDIN    TWITTER

Google's artificial intelligence (AI) ethics board, established when Google acquired London AI startup DeepMind in 2014, remains one of the biggest mysteries in tech, with both Google and DeepMind refusing to

### Ethics Advisory Panel

Each step forward for AI is a step into uncharted territory.

That's why we made it our mission to ask the complicated questions that don't have easy answers. And give birth to something no AI company had ever created before — the Ethics Advisory Panel. So when we build something, we aren't just asking if it's great for our customers. We're asking if it's great for humanity.

1. Ethics Isn't Just About Legal Risk
2. Internal vs. External Advisors: Pros And Cons
3. More than Lip-Service About Ethics
    -- Patrick Lin and Evan Selinger

Institute of Social Computing
소셜컴퓨팅연구소

28

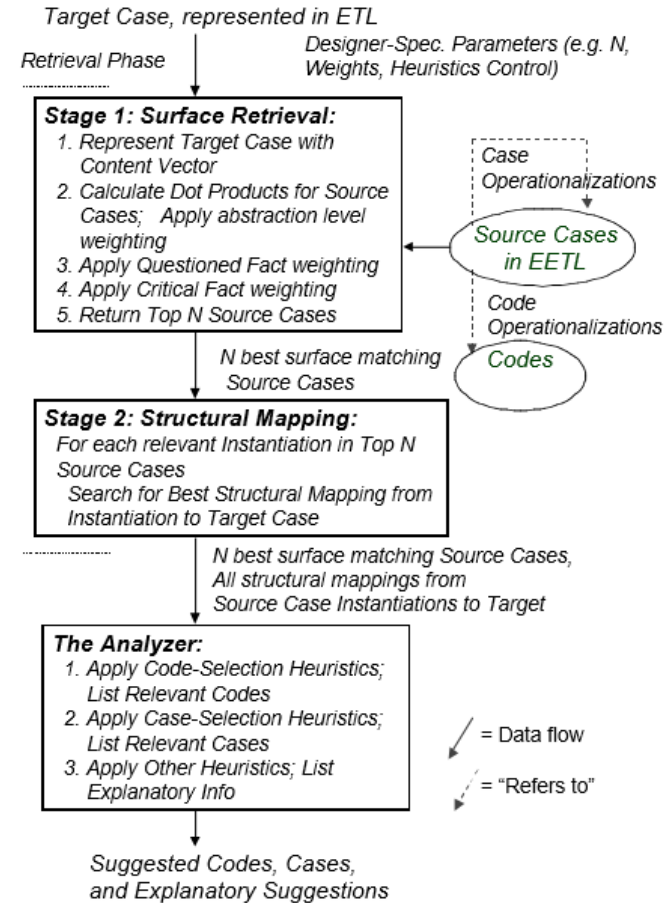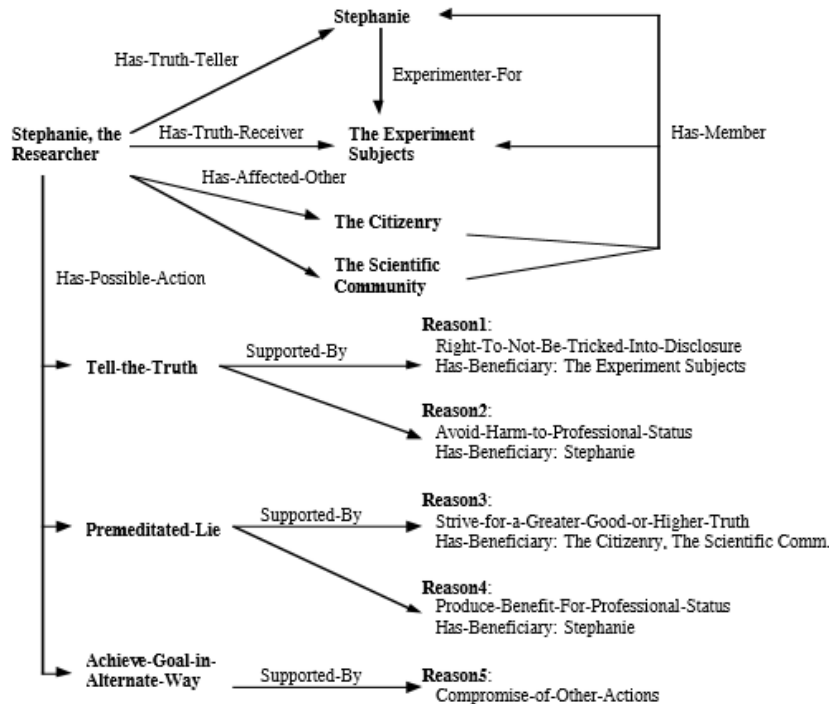# 모랄 코드는 구현이 가능할까?
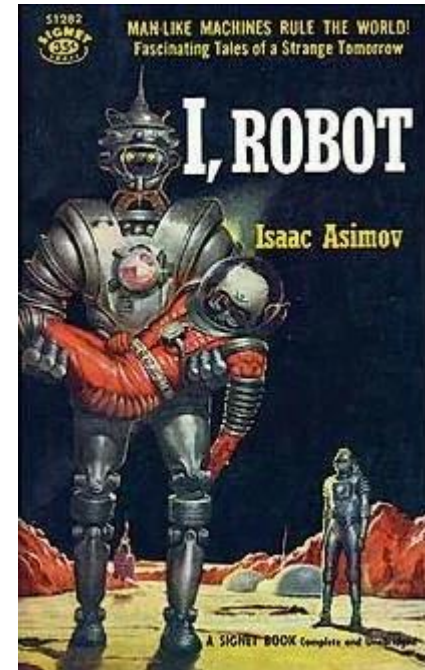


[출처: Reddit]

# Early Computational Models of Ethical Reasoning: Truth-Teller and SIROCCO



Bruce M. McLaren @ CMU

# 아시모프의 로봇 3원칙의 문제

1. **로봇은 인간에게 해를 가하거나 해를 당하는 상황에서 무시하면 안 된다**

   - '인간'을 어느 개념으로 정의할 것인지부터 쉬운 문제가 아니다. 인류는 한동안 다른 인종을 인간과 다르다고 분류한 적도 있고, 앞으로는 생물학적 특징만으로 인간을 정의하기 어려워질 수 있기 때문이다. 로봇에 이를 주입하는 것이 쉬운 일은 아니다.

   - '해를 가한다'는 것을 판단하려면 우선 그 행동의 결과가 누구에게 해가 될 수 있는지를 알아야 한다. 하지만 모든 상황을 정확히 판단한다는 것은 불가능한 일이다. 주변의 인간에게는 해를 가하지 않아도 지구 어딘가에 있는 다른 인간에게 해를 끼칠 수도 있는데 이를 계산할 방안이 없다.

2. **로봇은 1원칙에 어긋나지 않는 한, 인간의 명령에 복종해야 한다**

   - 인간의 발언 중 어디까지가 명령임을 알아내는 것 자체가 쉬운 일이 아니다. 제1원칙과 제2원칙을 따라서 행동한 것이 특정인을 구했지만, 그 결과로 인류에게 엄청난 파국을 일으킨다면 어떻게 할 것인가?

3. **로봇은 1, 2원칙에 어긋나지 않는 한, 자신을 지켜야 한다.**

▸ o원칙 (1985): **로봇은 인류에게 해를 가하거나, 또는 해를 당하는 상황을 무시해서는 안 된다.**

➔ 의무론적이고 하향식 방법의 문제는 규칙을 완벽하게 따를 때도 여전히 끔찍한 결과를 일으킬 수 있다

Institute of
Social Computing
소셜컴퓨팅연구소

# Value learning

**Question: How to specify complicated human values and ethics to AI systems?**

## Value learning by human feedback

**Stuart Russell:** Teach the agent by demonstrating human actions (cooperative inverse reinforcement learning).

**Owain Evans:** Human actions are often inconsistent and suboptimal. Modify inverse reinforcement learning to account for human biases.

**Paul Christiano:** Use semi-supervised learning to decrease reliance on human feedback (scalable AI control).

Viktoriya Krakovna, FLI / DeepMind @BAI2017

Institute of
Social Computing
소셜컴퓨팅연구소

# Value Learning

## Value learning by building in morality

**Francesca Rossi:** Specify ethical laws through <u>constraints</u>.

**Vincent Conitzer:** <u>Find patterns</u> in human ethical decisions, and build those features into AI systems.

**Adrian Weller:** Can we make human moral concepts more <u>precise and consistent</u>?

Viktoriya Krakovna, FLI / DeepMind @BAI2017

33

# 상향식 접근: Learn by Observation
## "How to Stop Your Robot Cooking Your Cat"



Stuart Russell is a professor of computer science at the University of California, Berkeley, and an expert on artificial intelligence.

COMMENTARY

**SHOULD WE FEAR SUPERSMART ROBOTS?**

If we're not careful, we could find ourselves at odds with determined, intelligent machines whose objectives conflict with our own

*By Stuart Russell*

• • • • • • • • • • • • • • • • • • • • • • • • • • • •

T IS HARD TO ESCAPE THE NAGGING SUSPICION THAT CREATING MACHINES smarter than ourselves *might* be a problem. After all, if gorillas had accidentally created humans way back when, the now endangered primates probably would be wishing they had not done so. But *why*, specifically, is advanced artificial intelligence a problem?

Hollywood's theory that spontaneously evil machine consciousness will drive armies of killer robots is just silly. The real problem relates to the possibility that AI may become incredibly good at achieving something other than what we really want. In 1960 legendary mathematician Norbert Wiener, who founded the field of cybernetics, put it this way: "If we use, to achieve our purposes, a mechanical agency with whose operation we cannot efficiently interfere…, we had better be quite sure that the purpose put into the machine is the purpose which we really desire."

A machine with a specific purpose has another property, one that we usually associate with living things: a wish to preserve its own existence. For the machine, this trait is not innate, nor is it something introduced by humans; it is a logical consequence of the simple fact that the machine cannot achieve its original purpose if it is dead. So if we send out a robot with the sole directive of fetching coffee, it will have a strong incentive to ensure success by disabling its own off switch or even exterminating anyone who might interfere with its mission. If we are not careful, then, we could face a kind of global chess match against very determined, superintelligent machines whose objectives conflict with our own, with the real world as the chessboard.

The prospect of entering into and losing such a match should concentrate the minds of computer scientists. Some researchers argue that we can seal the machines inside a kind of fire wall, using them to answer difficult questions but never allowing them to affect the real world. (Of course, this means giving up on superintelligent robots!) Unfortunately, that plan seems unlikely to work: we have yet to invent a fire wall that is secure against ordinary humans, let alone superintelligent machines.

Can we instead tackle Wiener's warning head-on? Can we de-

58 Scientific American, June 2016

---

- The machine's purpose must be to maximize the realization of human values. In particular, it has no purpose of its own and no innate desire to protect itself.

- The machine must be initially uncertain about what those human values are. The machine may learn more about human values as it goes along, of course, but it may never achieve complete certainty.

- The machine must be able to learn about human values by observing the choices that we humans make.

  - Inverse reinforcement learning (IRL), concerned with learning the values of some by observing its behavior. By watching a typical human's morning routine, the robot learns about the value of coffee to humans.

- 상향식 접근은 인간 행동의 목적과 결과, 영향, 행동이 윤리적 기반을 갖는 것임을 판단하는 능력을 갖추고, 이를 다시 내부의 코드로 만들어가야 하는 어려움을 갖고 있다.

- 인간 가치를 표현하고 이를 각 인간의 배경에 따라 윤리, 법, 도덕으로 인지할 수 있는 능력을 포함해야 하므로 아직 이런 방향의 연구는 초기 단계라고 생각한다. 더군다나 이런 인간 행동이 윤리적 의미를 갖게 될 때에는 다양한 인간 감성 표현과 감정을 이해해야 하는데, 이는 아직 우리가 알고 있는 인공지능 기술에서도 매우 어려운 분야이다.

Institute of Social Computing
소셜컴퓨팅연구소

# 상황 인지의 어려움

# 감정을 어떻게 인식할 것인가?
## Affective Computing

▸ The field of study concerned with understanding, recognizing, and utilizing human emotions and other affective phenomena in the design of technological systems - IEEE

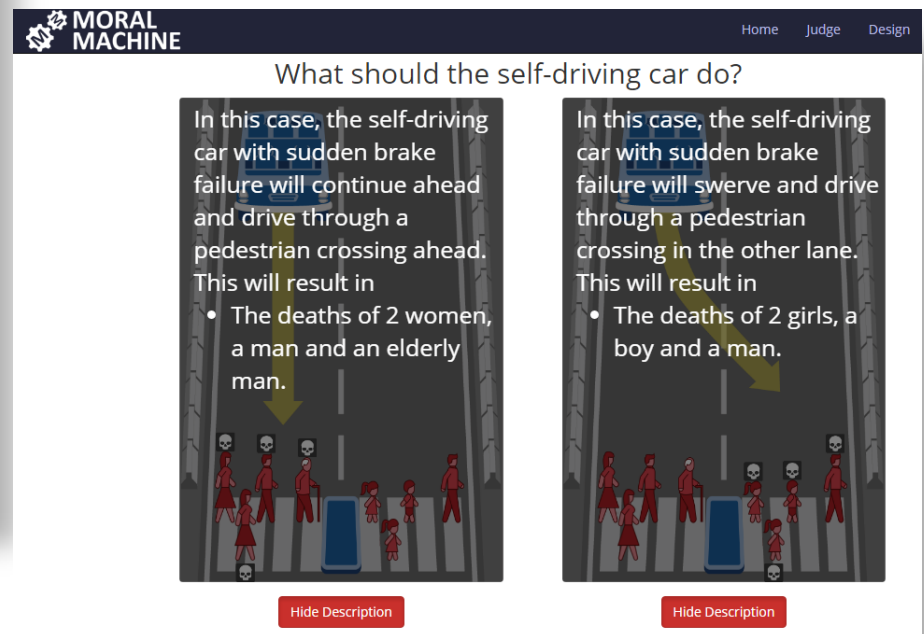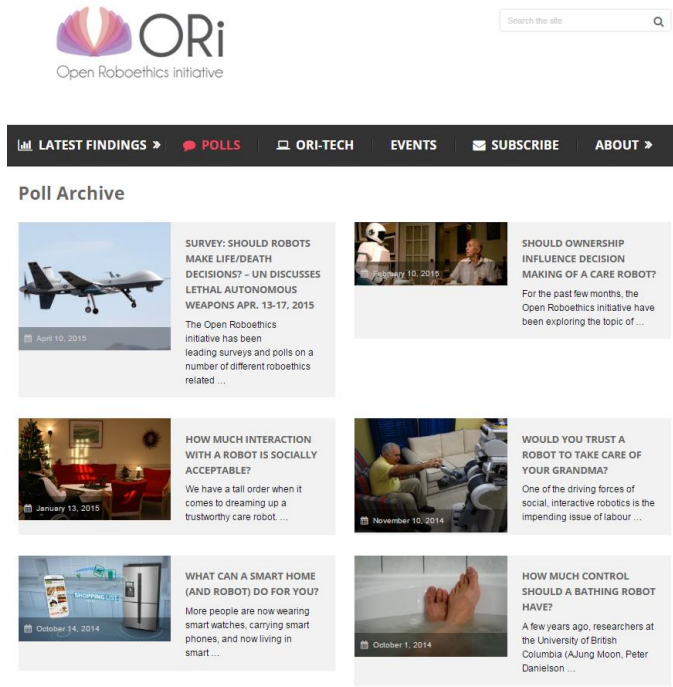# Using Stories to Teach Human Values to Artificial Agents

- Mark Riedl and Brent Harrison. 2016. Using Stories to Teach Human Values to Artificial Agents. In Proceedings of the 2nd International Workshop on AI, Ethics and Society.

- 로봇이나 인공 행위자(에이전트)가 이야기를 읽고, 각 사건의 바람직한 결과를 학습해, 인간 사회에서 성공적인 행동을 이해하도록 훈련하는 시스템

# 사람들의 참여에 의한 윤리 판단 데이터 수집

Providing a platform for 1) building a crowd-sourced picture of human opinion on how machines should make decisions when faced with moral dilemmas, and 2) crowd-sourcing assembly and discussion of potential scenarios of moral consequence





http://moralmachine.mit.edu/

# AI Safety Research Teams

Current AI safety research teams

Academia:

CENTRE FOR THE STUDY OF EXISTENTIAL RISK
UNIVERSITY OF CAMBRIDGE

CFI LEVERHULME CENTRE FOR THE FUTURE OF INTELLIGENCE

Future of Humanity Institute
UNIVERSITY OF OXFORD

UC Berkeley Center for Human-Compatible AI

**FLI grantees**

Independent:

**MIRI**
MACHINE INTELLIGENCE
RESEARCH INSTITUTE

Industry: DeepMind OpenAI

Institute of Social Computing
소셜컴퓨팅연구소

# Challenges and Further Researches

▶ Transparency

▶ Explainability

▶ Computational model of ethics

▶ How Do We Align Artificial Intelligence with Human Values?

◦ Understanding what "we" want

◦ Humanity do not agree on common values, and even parts we do agree on change with time

▶ Collaboration between AI researchers and Ethicists

Institute of
Social Computing
소셜컴퓨팅연구소

# 감사합니다
## (Meet me at
## 페이스북: facebook.com/stevehan)

# 소셜 로봇에 대한 법적 보호와 안전성 요구 문제



[출처: Kate Darling, MIT]

# Civil Law Rules on Robotics: Draft Report of EU Parliament

*Committee on Legal Affairs*

2015/2103(INL)

31.5.2016

## DRAFT REPORT

with recommendations to the Commission on Civil Law Rules on Robotics (2015/2103(INL))

Committee on Legal Affairs

Rapporteur: Mady Delvaux

(Initiative – Rule 46 of the Rules of Procedure)

---

**MOTION FOR A EUROPEAN PARLIAMENT RESOLUTION**

with recommendations to the Commission on Civil Law Rules on Robotics (2015/2103(INL))

*The European Parliament*,

– having regard to Article 225 of the Treaty on the Functioning of the European Union,

– having regard to Rules 46 and 52 of its Rules of Procedure,

– having regard to the report of the Committee on Legal Affairs and the opinions of the Committee on Employment and Social Affairs, the Committee on the Environment, Public Health and Food Safety, the Committee on Industry, Research and Energy and the Committee on the Internal Market and Consumer Protection (A8-0000/2016),

**Introduction**

A. whereas from Mary Shelley's Frankenstein's Monster to the classical myth of Pygmalion, through the story of Prague's Golem to the robot of Karel Čapek, who coined the word, people have fantasised about the possibility of building intelligent machines, more often than not androids with human features;

B. whereas now that humankind stands on the threshold of an era when ever more sophisticated robots, bots, androids and other manifestations of artificial intelligence ("AI") seem poised to unleash a new industrial revolution, which is likely to leave no stratum of society untouched, it is vitally important for the legislature to consider all its implications;

C. whereas between 2010 and 2014 the average increase in sales of robots stood at 17% per year and in 2014 sales rose by 29%, the highest year-on-year increase ever, with automotive parts suppliers and the electrical/electronics industry being the main drivers of the growth; whereas annual patent filings for robotics technology have tripled over the last decade;

D. whereas in the short to medium term robotics and AI promise to bring benefits of efficiency and savings, not only in production and commerce, but also in areas such as transport, medical care, education and farming, while making it possible to avoid exposing humans to dangerous conditions, such as those faced when cleaning up toxically polluted sites; whereas in the longer term there is potential for virtually unbounded prosperity;

E. whereas at the same time the development of robotics and AI may result in a large part of the work now done by humans being taken over by robots, so raising concerns about the future of employment and the viability of social security systems if the current basis of taxation is maintained, creating the potential for increased inequality in the distribution of wealth and influence;

F. whereas the causes for concern also include physical safety, for example when a robot's code proves fallible, and the potential consequences of system failure or hacking of
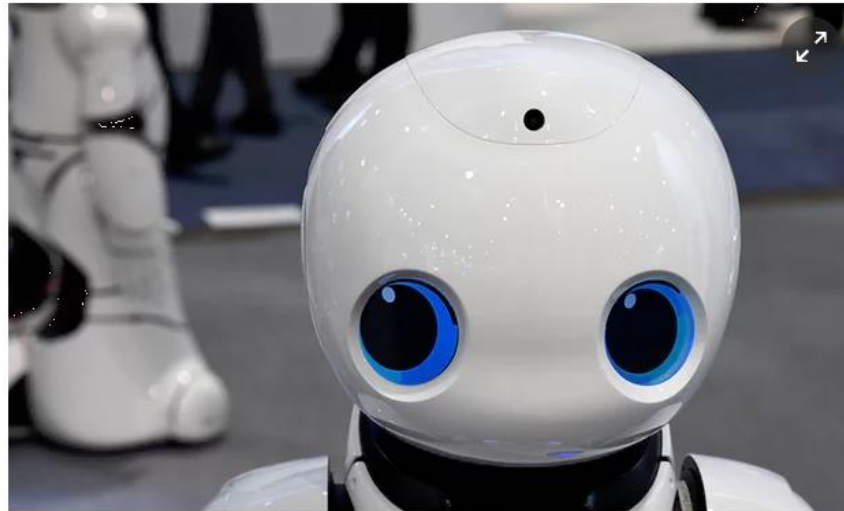
Institute of Social Computing
소셜컴퓨팅연구소

43

# Electronic Personhood?

## Give robots 'personhood' status, EU committee argues

Proposed rules for robots and AI in Europe include a push for a general basic income for humans, and 'human rights' for robots



ⓘ A Tanscorp UU smart robot is displayed at CES 2017 at the Sands Expo and Convention Center in Las Vegas. Photograph: Ethan Miller/Getty

The European parliament has urged the drafting of a set of regulations to govern the use and creation of robots and artificial intelligence, including a form of "electronic personhood" to ensure rights and responsibilities for the most capable AI.

In a 17-2 vote, with two abstentions, the parliament's legal affairs committee passed the report, which outlines one possible framework for regulation.

Institute of
Social Computing
소셜컴퓨팅연구소