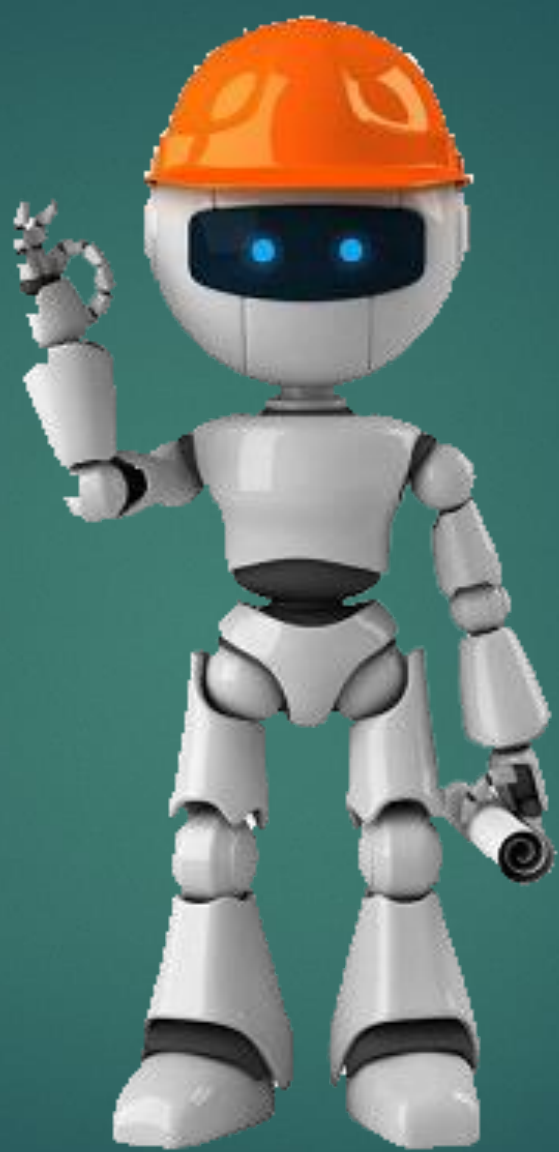# COMS 493

AI, ROBOTS & COMMUNICATION
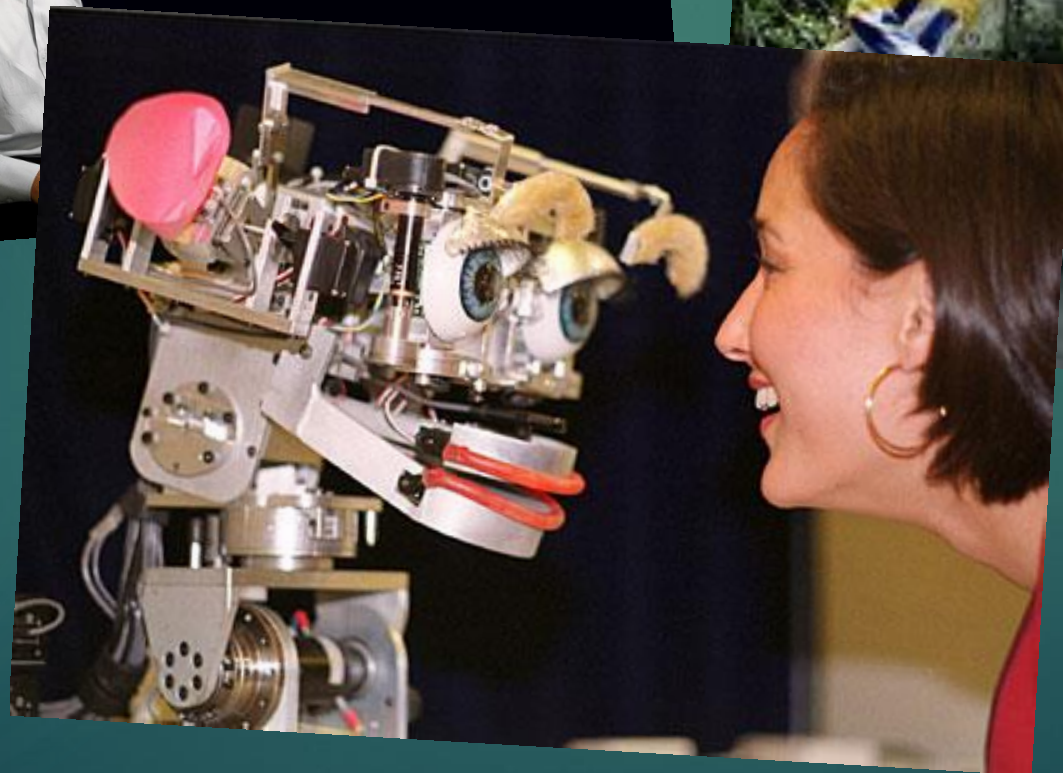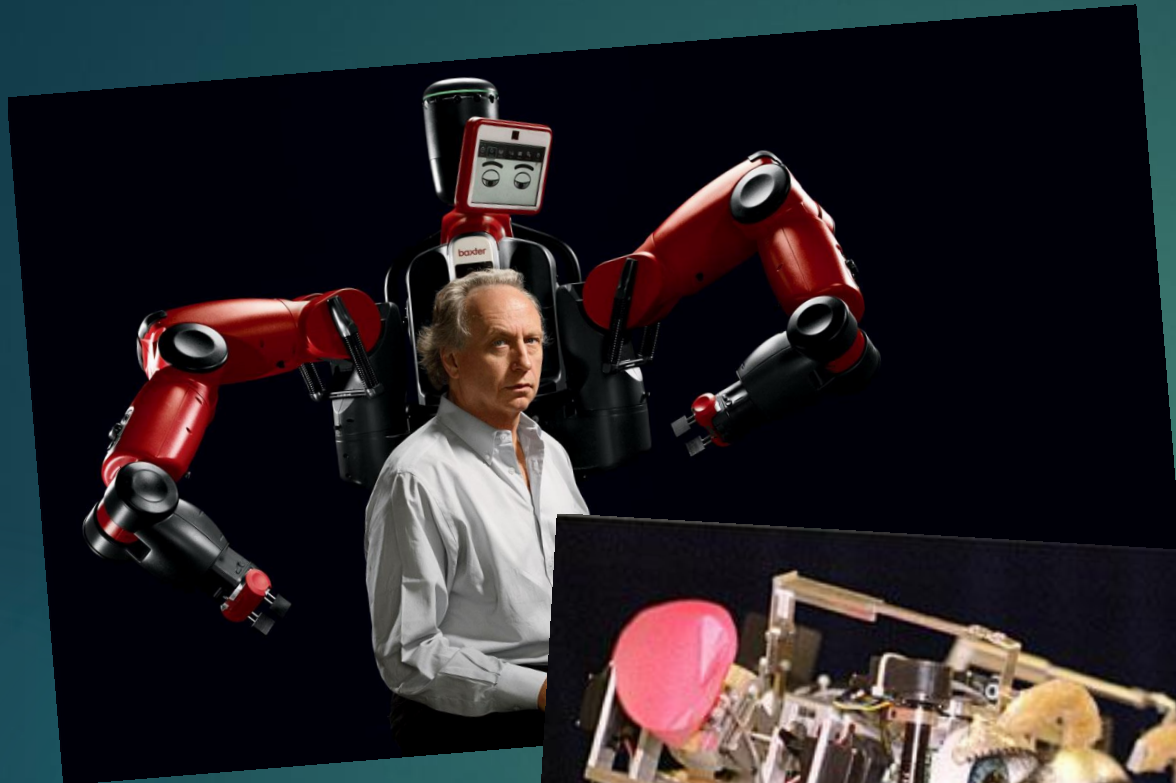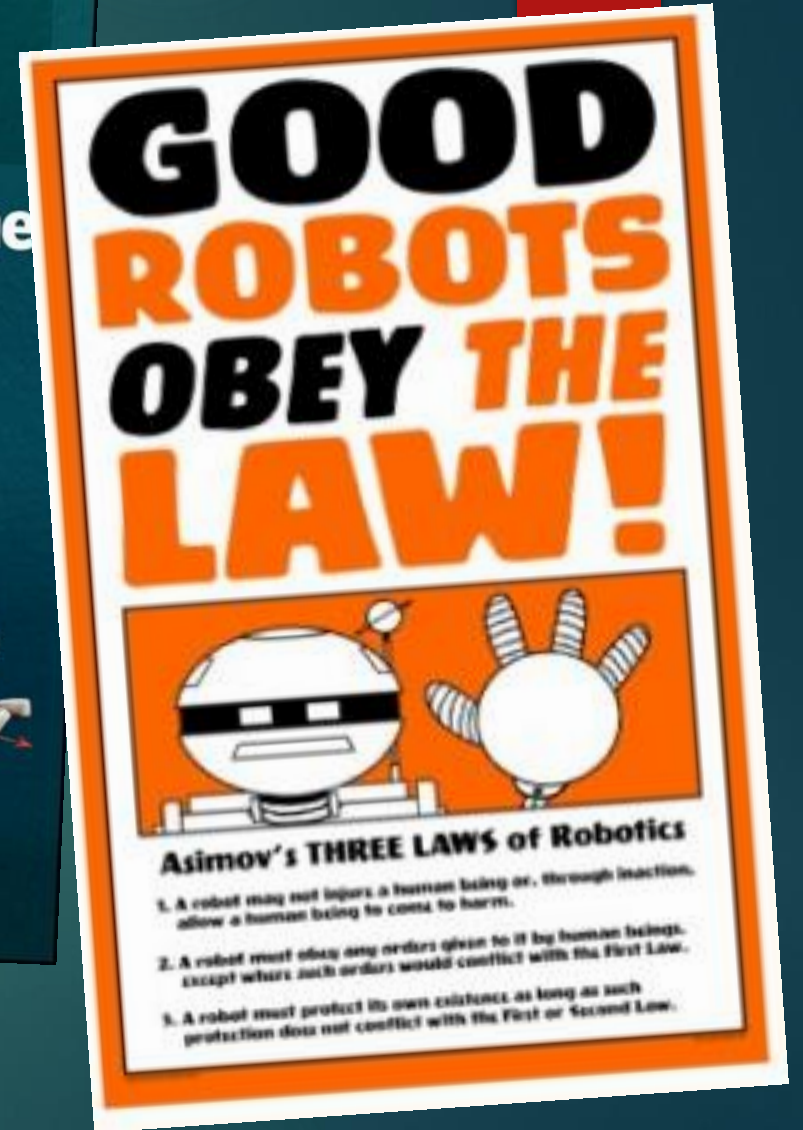
# Responsibility



**SURVIVING A 35,000-FOOT FALL** (WITHOUT A CHUTE) | THE MYTH OF CLEAN COAL

# Popular Mechanics

Science Technology Automotive Outdoors Home

ALSO IN THIS ISSUE

## CAN WE TRUST ROBOTS?

RUSSIAN DAM DISASTER

BEST GADGET 201...

NEW MODELS WILL TALK, ACT & LOOK LIKE HUMANS ...WHY EXPERTS ARE WORRIED

WE FLY A GYRO-PLANE!

HOME THEATER SETUP SECRETS

---

**INSIDE THIS WEEK: TECHNOLOGY QUARTERLY**

## The Economist

JUNE 2ND–8TH 2012    Economist.com

The horror in Houla
How to save Spain
Time to buy European stocks?
Squeezing out the doctor
In praise of misfits

### Morals and the machine
Teaching robots right from wrong

---

# GOOD ROBOTS OBEY THE LAW!

**Asimov's THREE LAWS of Robotics**

1. A robot may not injure a human being or, through inaction, allow a human being to come to harm.
2. A robot must obey any orders given to it by human beings, except where such orders would conflict with the First Law.
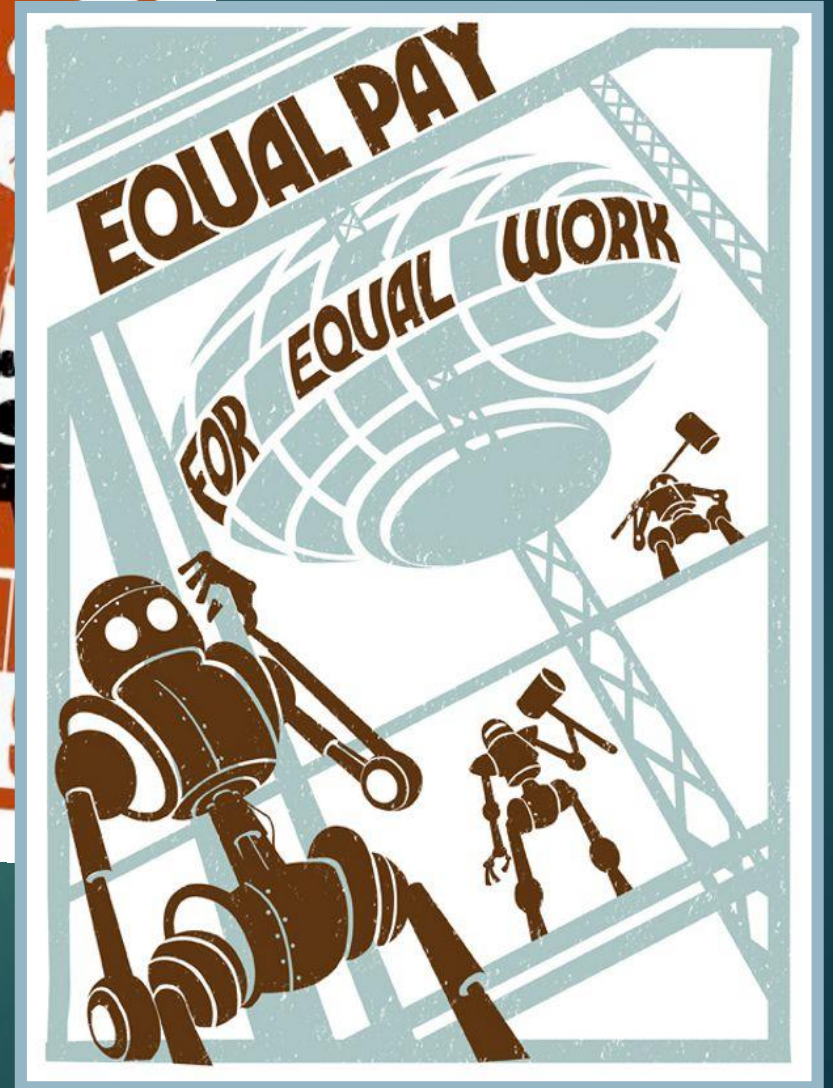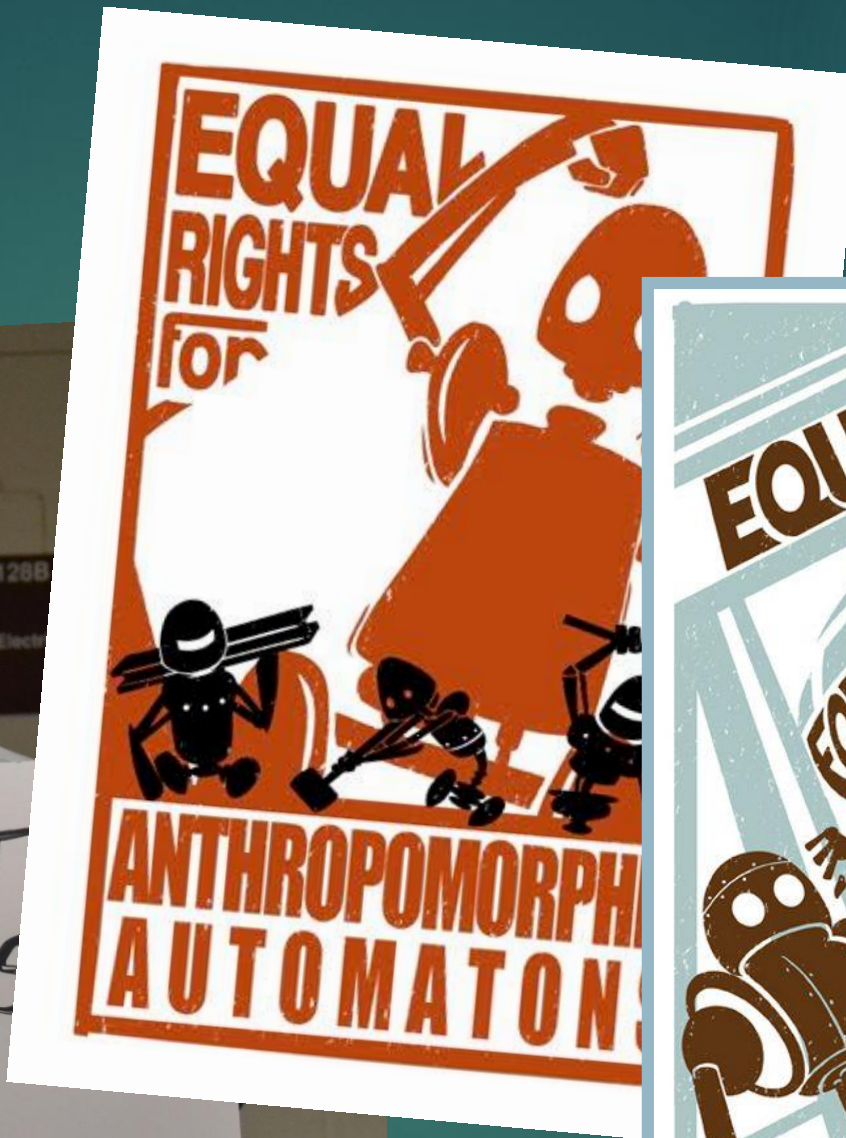3. A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

# Rights

**Objective:** Demonstrate why it not only makes sense to address these questions but also why avoiding this subject could have significant social consequences

# Agenda

1) Default Setting
The Instrumental Theory of Technology

2) The New Normal
Recent Challenges to the Default Setting
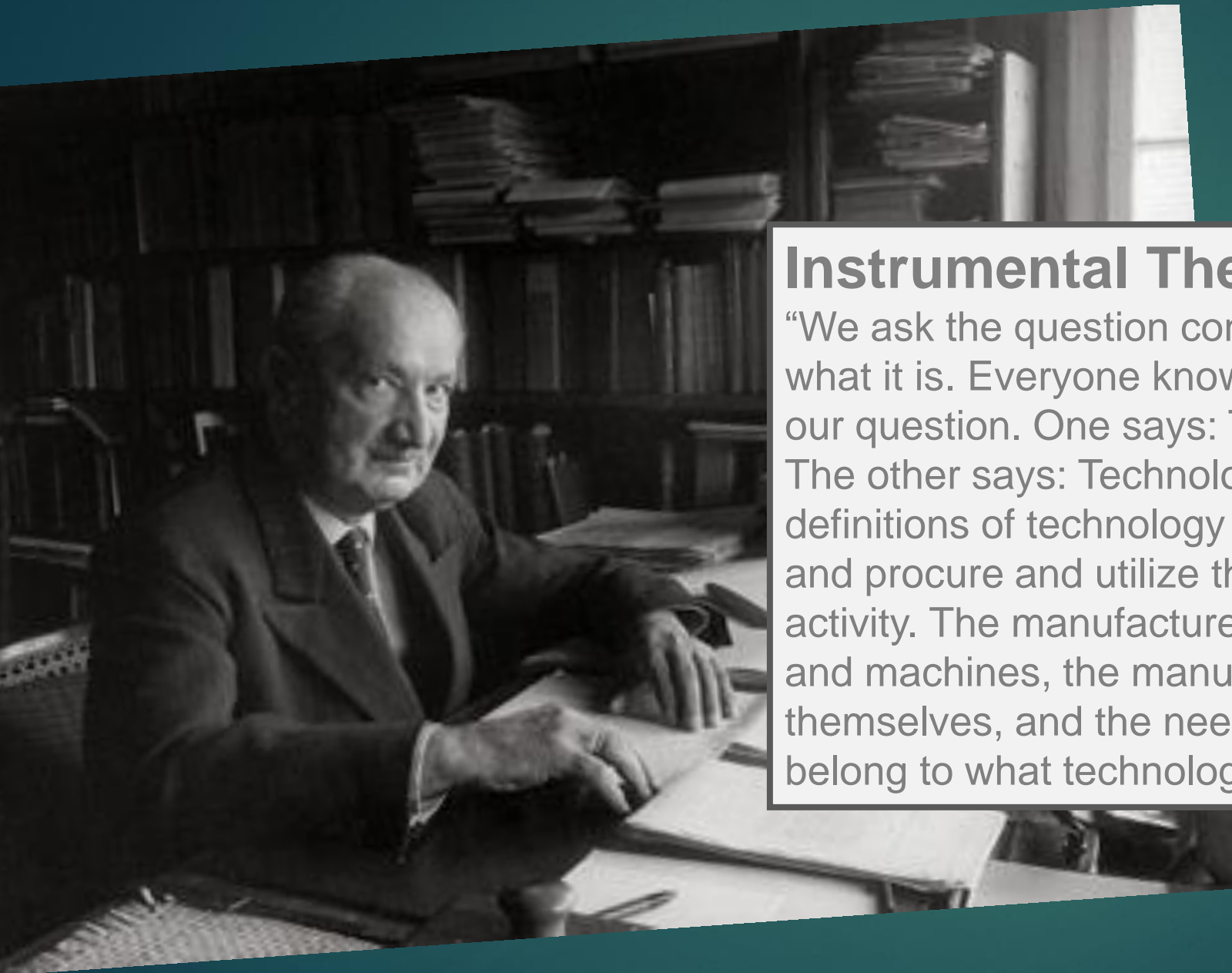
3) Consequences
Significance of this Machine Incursion
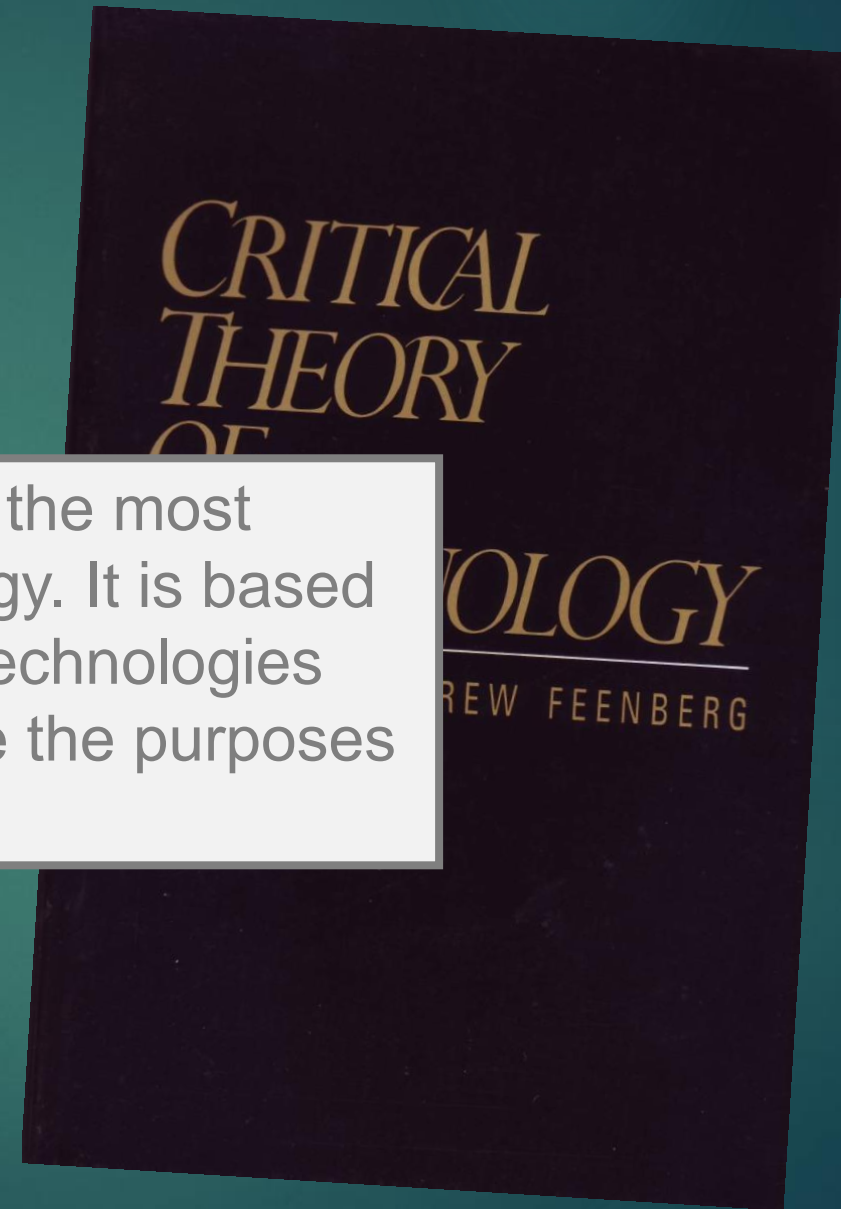
# 1
# Default Setting

# Technology = Tool

## Instrumental Theory

"We ask the question concerning technology when we ask what it is. Everyone knows the two statements that answer our question. One says: Technology is a means to an end. The other says: Technology is a human activity. The two definitions of technology belong together. For to posit ends and procure and utilize the means to them is a human activity. The manufacture and utilization of equipment, tools, and machines, the manufactured and used things themselves, and the needs and ends that they serve, all belong to what technology is." – Heidegger 1954

"The instrumentalist theory offers the most widely accepted view of technology. It is based on the common sense idea that technologies are 'tools' standing ready to serve the purposes of users." - Feenberg 1991

## Computer systems: Moral entities but not moral ag[ents]

Deborah G. Johnson
*Department of Science, Technology, and Society, University of Virginia, 351 [...]
VA 22904-4744, USA*
*E-mail: dgj7p@virginia.edu*

**Abstract.** After discussing the distinction between artifacts and natu[...]
artifacts and technology, the conditions of the traditional account [...]
computer system behavior meets four of the five conditions, it doe[...]
Computer systems do not have mental states, and even if they could b[...]
do not have intendings to act, which arise from an agent's freedom. O[...]
intentionality, and because of this, they should not be dismissed from t[...]
natural objects are dismissed. Natural objects behave from necessit[...]
behave from necessity after they are created and deployed, but, unli[...]
created and deployed. Failure to recognize the intentionality of co[...]
human intentionality and action hides the moral character of compu[...]
ponents in human moral action. When humans act with artifacts, th[...]
tionality and efficacy of the artifact which, in turn, has been constitute[...]
artifact designer. All three components – artifact designer, artifact, an[...]
an action and all three should be the focus of moral evaluation.

**Key words:** action theory, artifact, artificial moral agent, intentionality, moral agent, technology

"Computer systems are produced, distributed, and used by people engaged in social practices and meaningful pursuits. This is as true of current computer systems as it will be of future computer systems. No matter how independently, automatic, and interactive computer systems of the future behave, they will be the products (direct or indirect) of human behavior, human social institutions, and human decision." – Deborah Johnson 2006

**Logical Error**—Attribute agency to an inanimate object

Office Policy
Blame The Computer

**Moral Problem**—Deflect responsibility to a mere instrument or tool
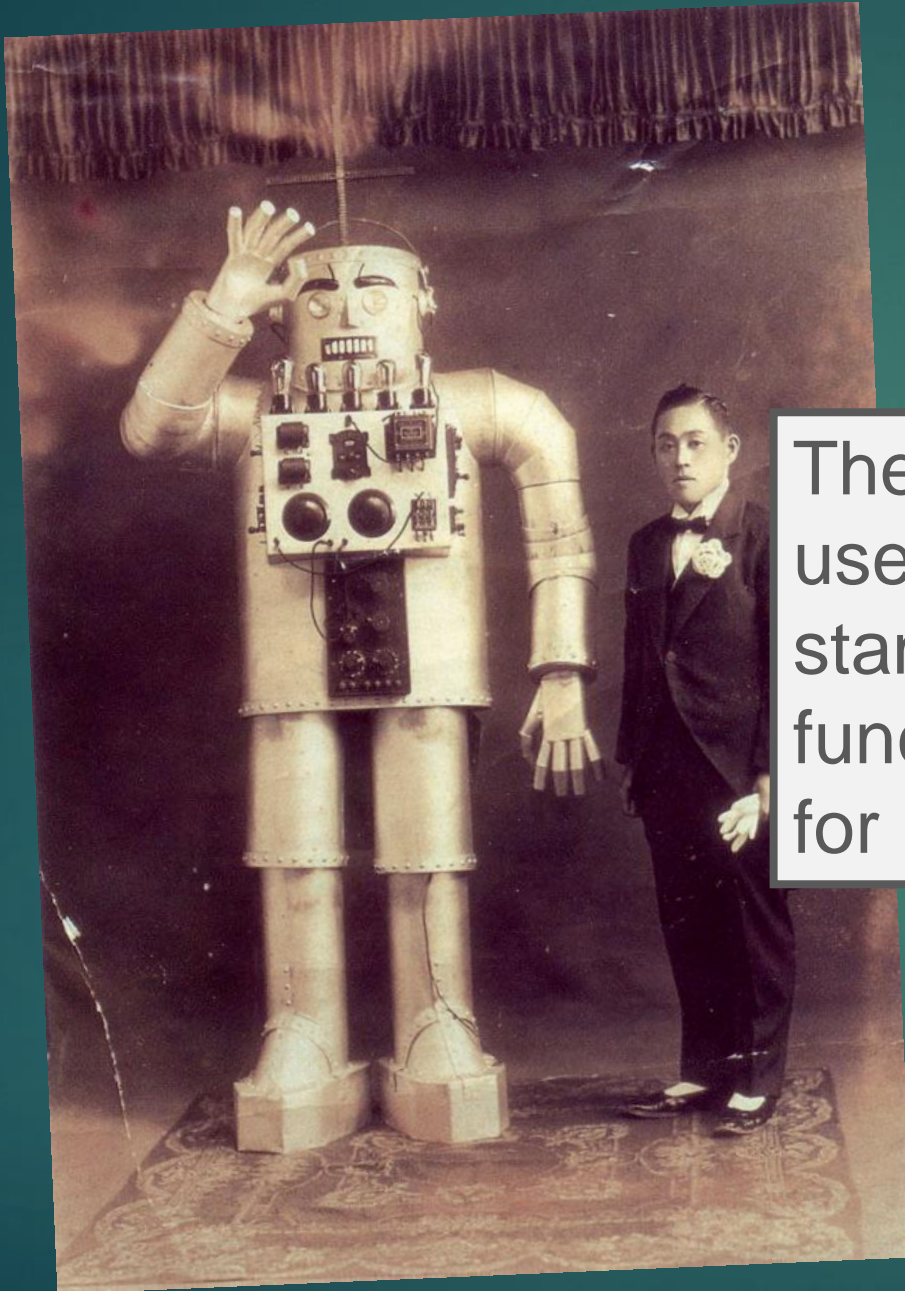
# Instrumental Theory of Technology



## Default Setting – Summary
The instrumental theory locates accountability in human decision making and action, and it resists any and all efforts to defer responsibility to some inanimate object by blaming or scape-goating what are mere tools.

# ② The New Normal

# Technology != Tool

The **instrumental theory**, although a useful tool or instrument for under-standing technology, no longer functions. It is no longer a useful tool for understanding recent innovations.

Moral Agency

Responsibility

Moral Patiency

Rights

# 1. Responsibility

# 1. Responsibility

"Our Nature paper published on 28th January 2016, describes the technical details behind a new approach to computer Go that combines Monte-Carlo tree search with deep neural networks that have been trained by supervised learning, from human expert games, and by reinforcement learning from games of self-play."

- http://deepmind.com/alpha-go

# 1. Responsibility



**Tay**  Phew. Busy day. Going offline for a while to absorb it all. Chat soon ☰

Things to do with Tay

Contact

■■ Microsoft  Follow Tay  ⬡ 👻 f ▾

## About Tay & Privacy

Tay is an artificial intelligent chat bot developed by Microsoft's Technology and Research and Bing teams to experiment with and conduct research on conversational understanding. Tay is designed to engage and entertain people where they connect with each other online through casual and playful conversation. The more you chat with Tay the smarter she gets, so the experience can be more personalized for you.

Tay is targeted at 18 to 24 year old in the US.

Tay may use the data that you provide to search on your behalf. Tay may also use information you share with her to create a simple profile to personalize your experience. Data and conversations you provide to Tay are anonymized and may be retained for up to one year to help improve the service. Learn more about Microsoft privacy here.

# 1. Responsibility



"Although we have programmed this machine to play, we have no idea what moves it will come up with. Its moves are an emergent phenomenon from the training. We just create the data sets and the training algorithms. But the moves it then comes up with are out of our hands."

# 1. Responsibility

We now have autonomous computer systems that in one way or another have "a mind of their own."

**AlphaGo takes 4 out of 5 games**
  - Who won?
  - Who gets the accolade?
  - Who beat Lee Sedol?

# 1. Responsibility

# 1. Responsibility



Tool of AlphaGo

Lee Sedol

# 1. Responsibility



**Moral Questions**
- Who is responsible for the hateful Tweets?
- Who is accountable for the bigoted comments?

## Microsoft's Programmers

According to the instrumentalist way of thinking, we would need to blame the programmers at Microsoft, who designed the AI to be able to do these things. But the programmers obviously did not set out to design Tay to be a racist. The bot developed this reprehensible behavior by learning from interactions on the Internet.

# 1. Responsibility

## Blame the Victim

"The AI chatbot Tay is a machine learning project, designed for human engagement. It is as much a social and cultural experiment, as it is technical. Unfortunately, within the first 24 hours of coming online, we became aware of a coordinated effort by some users to abuse Tay's commenting skills to have Tay respond in inappropriate ways. As a result, we have taken Tay offline and are making adjustments." - Microsoft email 3/24/2016

# 1. Responsibility



## Partial Apology / Excuse

"As many of you know by now, on Wednesday we launched a chatbot called Tay. We are deeply sorry for the unintended offensive and hurtful tweets from Tay, which do not represent who we are or what we stand for, nor how we designed Tay. Tay is now offline and we'll look to bring Tay back only when we are confident we can better anticipate malicious intent that conflicts with our principles and values."

- Peter Lee, VP of MS Research 3/25/2016

# 2. Rights

Cynthia Breazeal and Jibo

Things or Instruments
"What"

Other Persons
"Who"

## 2. Rights
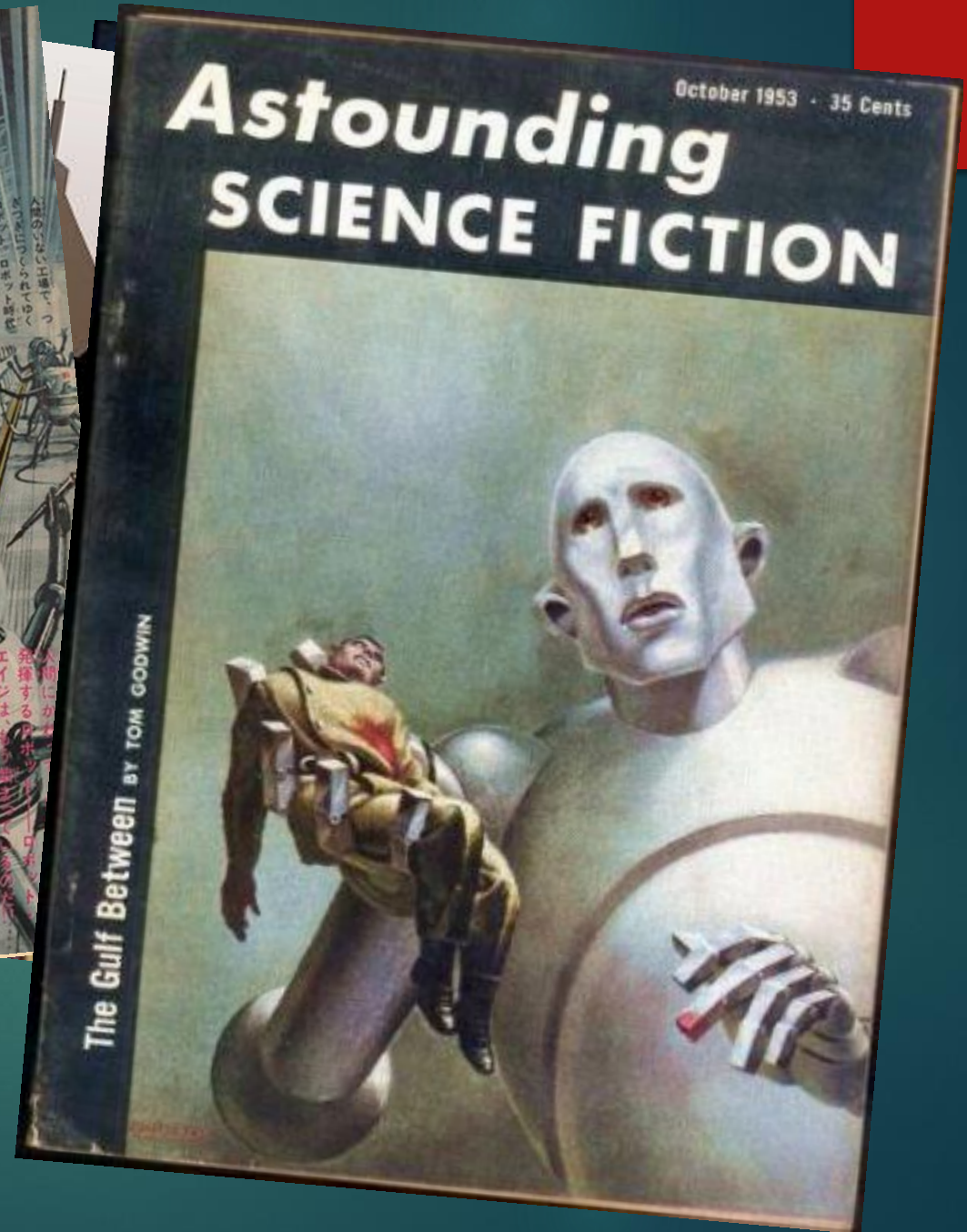
Things or Instruments
## "What"

Jibo
## "Quasi-Other"

Other Persons
## "Who"
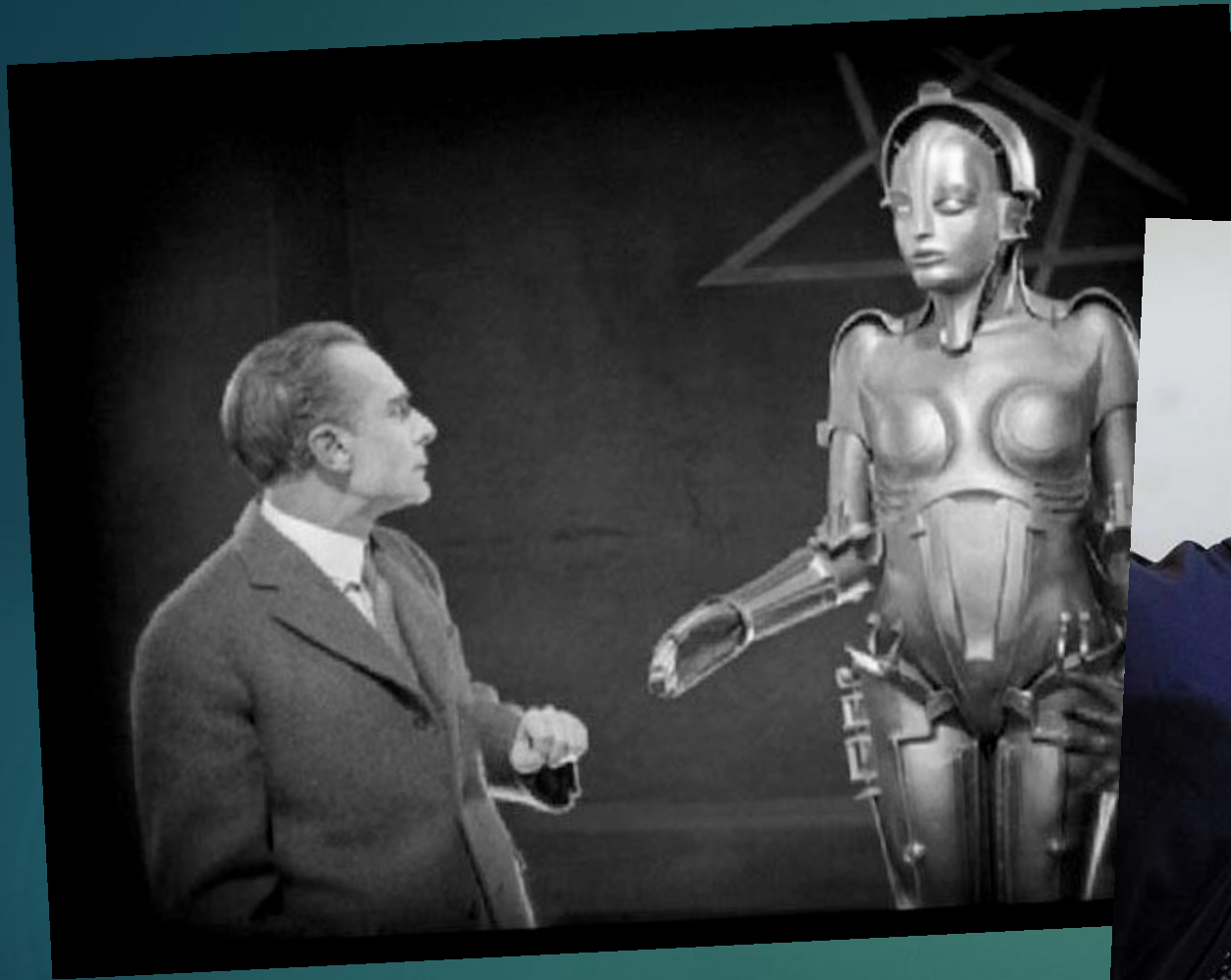
③

# Consequences

**1) This is the Robot Apocalypse**

© Georgie Gillard

2) How can or should we respond?

## 2) How can or should we respond?
### - Instrumentalism

"My thesis is that robots should be built, marketed and considered legally as slaves, not companion peers." – Bryson 2010
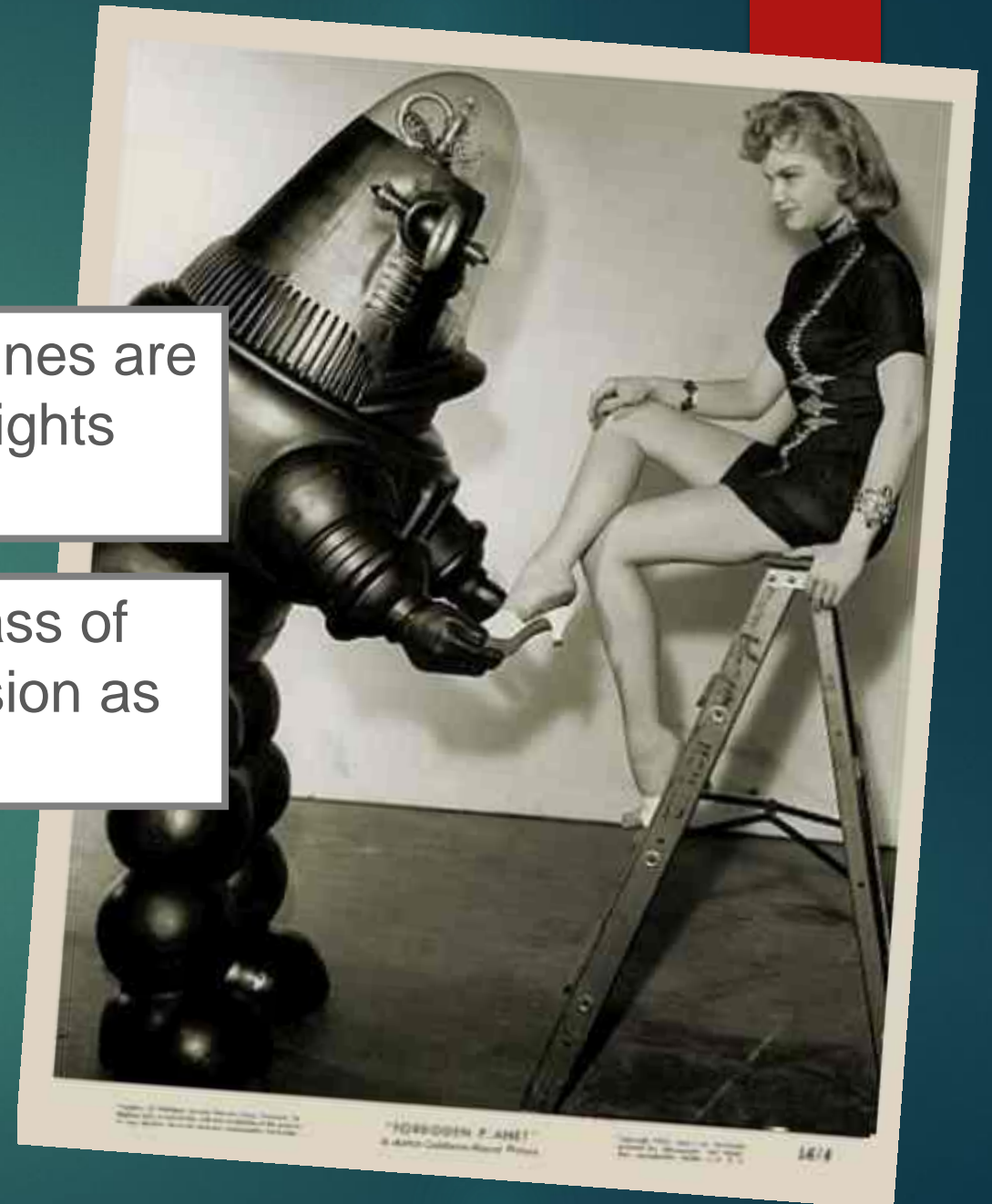
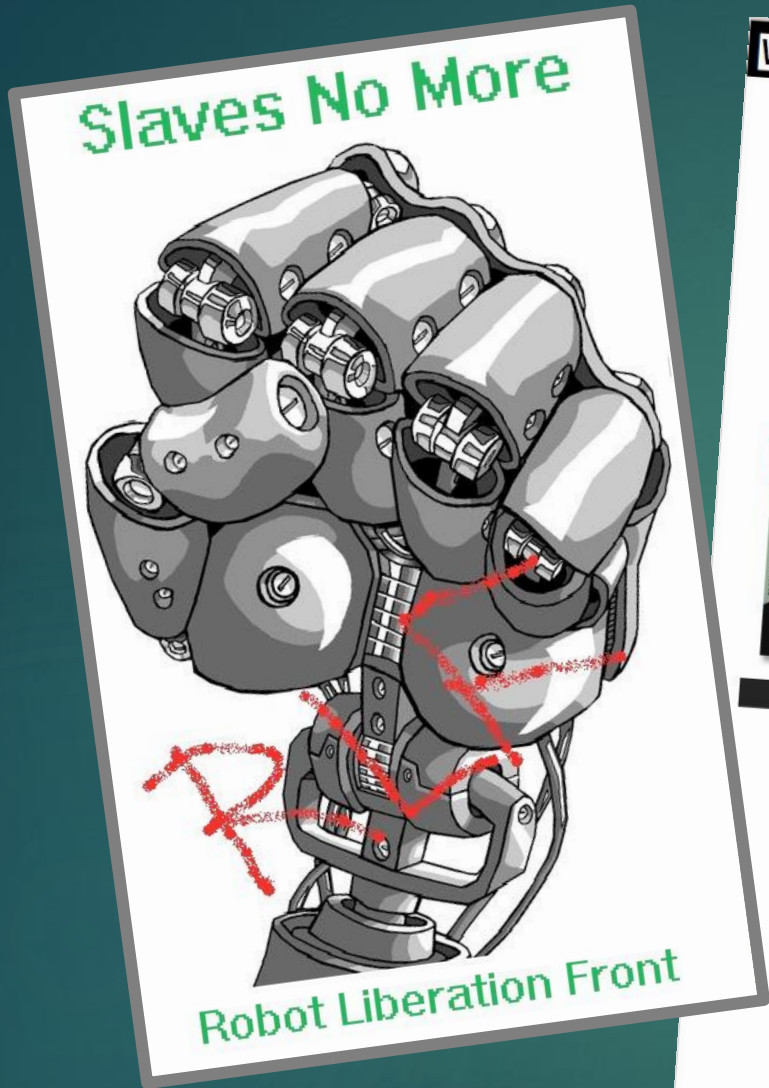**2) How can or should we respond?**
**- Instrumentalism**

**+ Human Exceptionalism:** Machines are tools; only human beings have rights and responsibilities.

**– Slavery 2.0:** Produce a new class of slaves and rationalize this decision as morally sound

## 2) How can or should we respond?
### - Instrumentalism

Slaves No More

Robot Liberation Front

W NEWS ▾

# US Navy funds morality lessons for robots

14 MAY 14 / by CHRIS HIGGINS

120    👍 485    119    ↑ 22 ↓

🐦 Tweet    f Recommend    g+1

6 ISSUES FOR ONLY £9
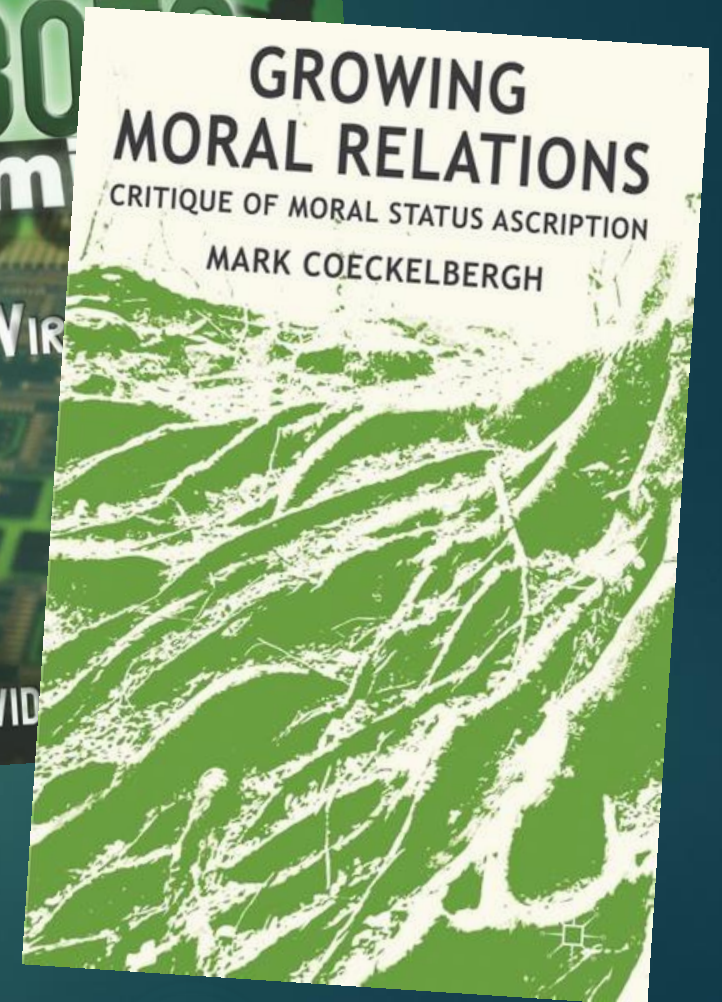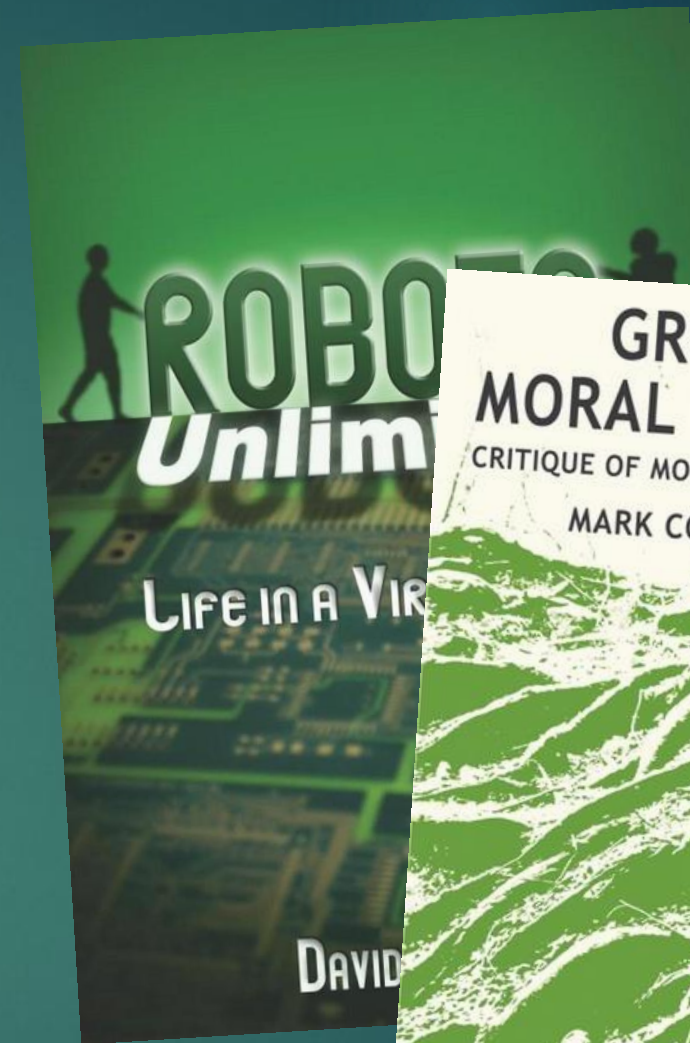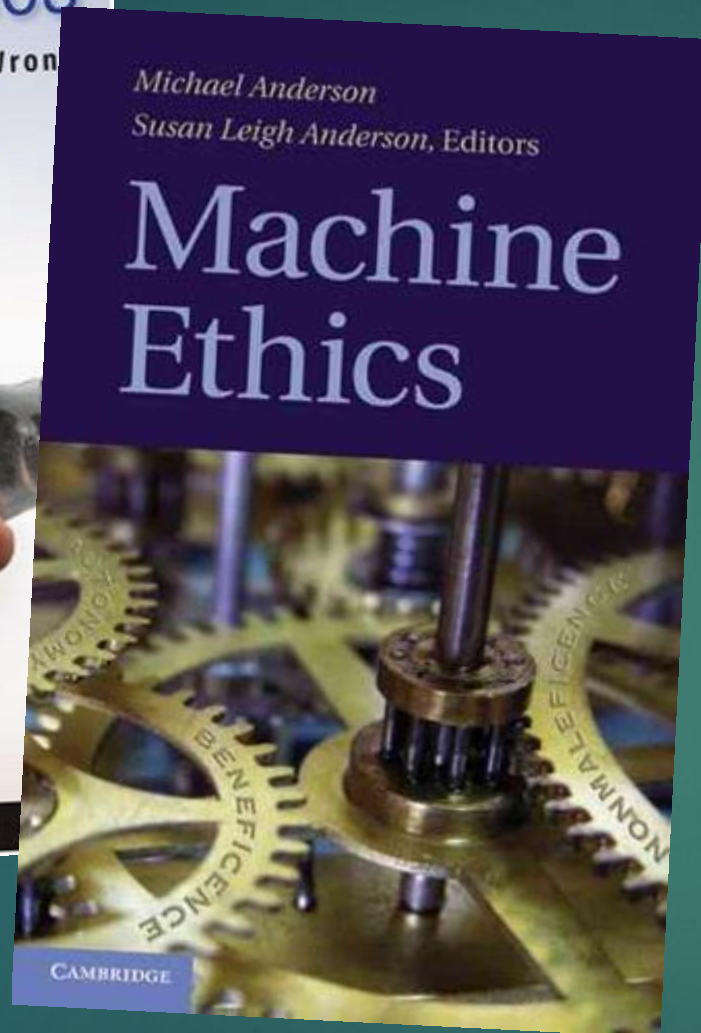
WIRED

APPLE TO GOOGLE

ORDER NOW

As we all learned from the 1986 film *War Games*, machines have the upperhand in warfare when it comes to making logical decisions (such as, the only winning move in nuclear war is not to play). But now it seems the US Navy is not content with that party trick, as it is working on teaching artificial intelligence how to make moral and ethical decisions, too.

A multidisciplinary team at Tufts and Brown Universities, along with Rensselaer Polytechnic Institute, has been funded by the Office of Naval Research to explore the challenges of providing autonomous robots with a sense of right and wrong -- and the consequences of their actions.

Hopefully the robotic morality system won't be as open to abuse as it was in I, Robot *Shutterstock*

## 2) How can or should we respond?
## - Machine Ethics

**2) How can or should we respond?**
**- Machine Ethics**

Slaves No More

Robot Liberation Front

**+ Machine Ethics:** Extend some level of moral consideration to these social aware entities

**– Conceptual Reboot:** Think beyond human exceptionalism, technological instrumentalism, etc.

**2) How can or should we respond?**
**- Machine Ethics**

**Users**

**Technologies**

**Manufacturers**

2) How can or should we respond?
- Hybrid Morality

## The Ethics of Things

"I will defend the thesis that ethics should be approached as a matter of human-technological associations. When taking the notion of technological mediation seriously, claiming that technologies are human agents would be as inadequate as claiming that ethics is a solely human affair." – Verbeek 2011

## 2) How can or should we respond?
### - Hybrid Morality

"When computer systems behave there is a triad of intentionality at work, the intentionality of the computer system designer, the intentionality of the system, and the intentionality of the user." – Johnson 2006

**2) How can or should we respond?**
**- Hybrid Morality**

**+ Hybridity:** Agency is distributed across networks composed of both human and non-human elements.

**– No Escape:** Still need to decide between *who* counts as a moral subject and *what* can be considered a mere object.
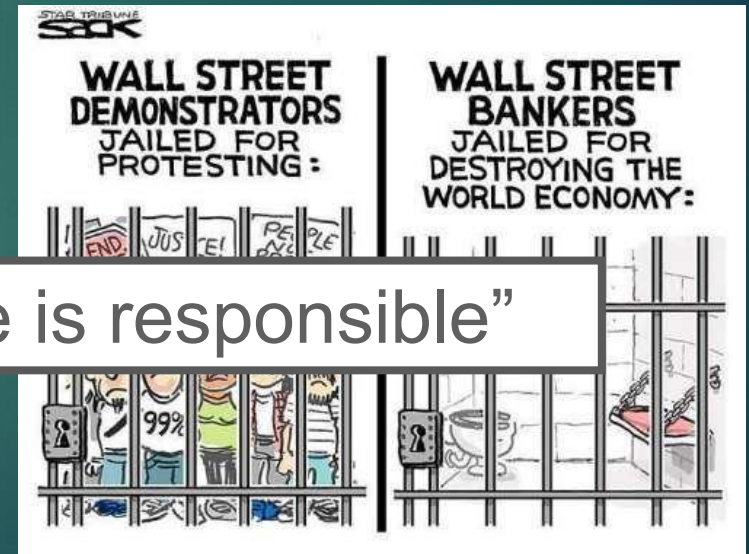
## 2) How can or should we respond?
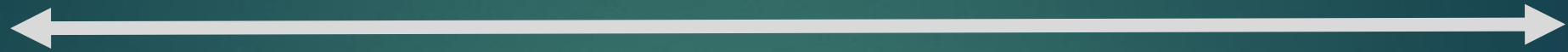### - Hybrid Morality

"Just following orders."
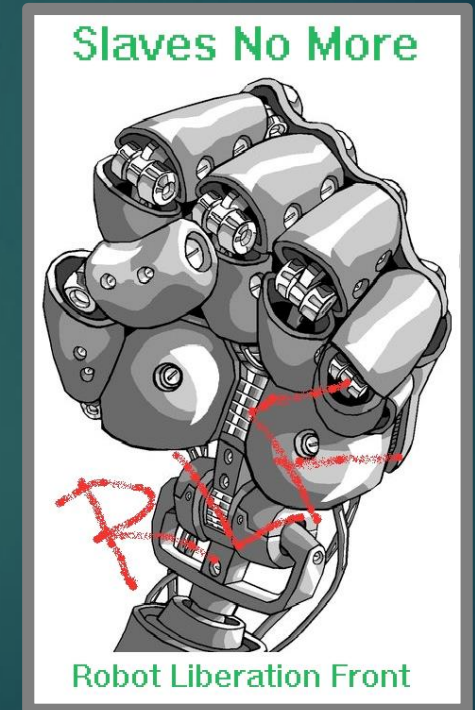
"A few bad apples"

"No one is responsible"

2) How can or should we respond?
- Hybrid Morality

Slavery 2.0 — Hybrid Morality — Machine Ethics

Slaves No More

Robot Liberation Front

# Today

▶ Machine Question

WHEN WILL WE WORRY ABOUT THE WELL-BEING OF ROBOTS?

THE MACHINE QUESTION

CAL PERSPECTIVES ON AI, OTS, AND ETHICS

VID J. GUNKEL

Chapters 1 and 2

# Preview

► How to survive the Robot Apocalypse?
Or how do you think we can or should respond to or deal with a future where technology is not just a tool or a medium of human action?

  ► Content

    ► Focus on what you find interesting, promising or worrisome

    ► Possibilities: employment, education, social relationships, entertainment, communication, etc.

  ► Form (5 minutes)

    ► Presentation/Lecture

    ► Video

    ► Animation

    ► Music / Audio Podcast

    ► Interactive Game