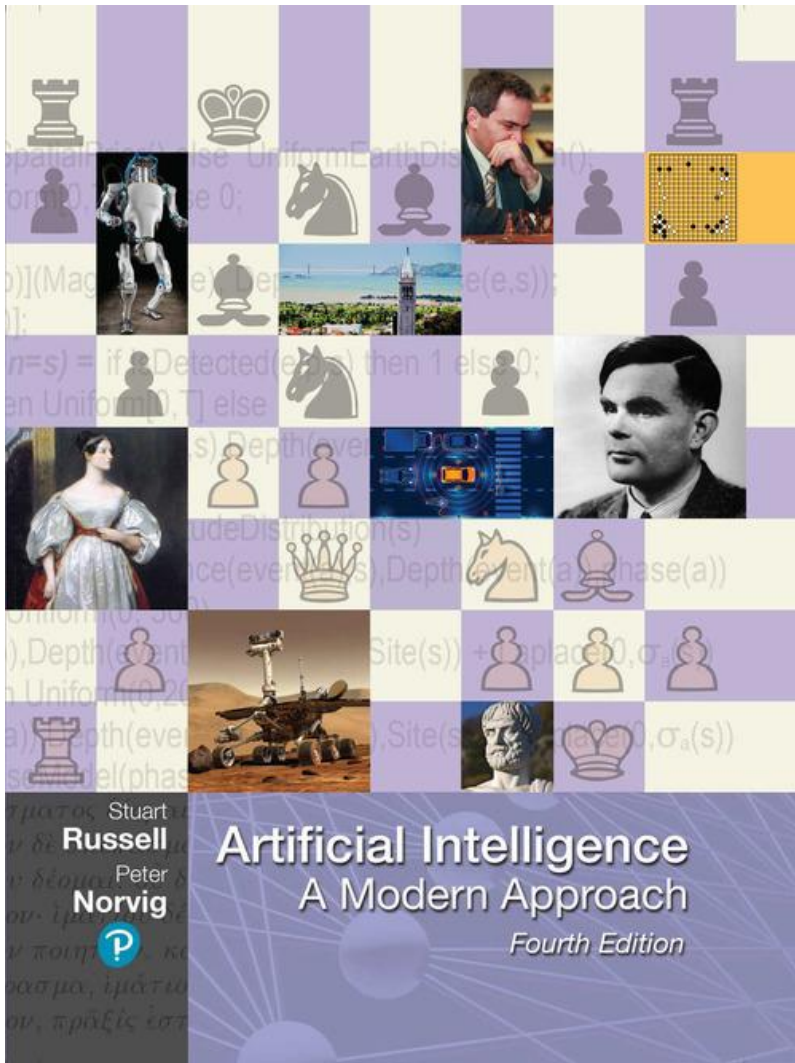# Artificial Intelligence Fundamentals

## 2024-2025



*"Success in creating AI would be the biggest event in human history. Unfortunately, it might also be the last, unless we learn how to avoid the risks."*

*- Stephen Hawking*

## AIMA Chapters 28-29

Philosophy, Ethics, Safety and Future of AI

# Outline

♦ The Limits of AI

♦ Can Machines Really Think?

♦ The Ethics of AI

♦ Future AI

# Course vs Curriculum vs AIMA

ML, ISPR, Comp. Neuro

Robotics, ISPR, NLP, etc.

Social and ethical issues in computer technology

Game Theory, Parallel & Distrib Sys.

Exercises (website)
Figures (pdf)
Code (website); Pseudocode (pdf)
Covers: US, Global

**ARTIFICIAL INTELLIGENCE: A MODERN APPROACH, 4TH US ED. (BERKELEY.EDU)**

# The Limits of AI

Philosopher *John Searle* (1980):

- **Weak AI**: the idea that machines could act as if they were intelligent
- **Strong AI**: the assertion that machines that do so are *actually consciously* thinking (not just *simulating thinking*)

The argument from informality
Turing's "*argument from informality of behavior*" says that human behavior is far **too complex** to be captured by any formal set of rules

**Good Old-Fashioned AI (GOFAI)**

- Simplest logical agent design
- *Qualification problem*: difficult to capture every contingency of appropriate behavior in a set of necessary and sufficient logical rules
- *Hubert Dreyfus*'s *"What Computers Can't Do"* strongest arguments is for situated agents rather than disembodied *logical inference engines*
- **Embodied cognition** approach claims that it makes no sense to consider the brain separately
  - Cognition takes place within a body, which is embedded in an environment

# The Limits of AI

The argument from disability

The "*argument from disability*" makes the claim that "*a machine can never do X.*" Turing's lists of *X*:

*Be kind, resourceful, beautiful, friendly, have initiative, have a sense of humor, tell right from wrong, make mistakes, fall in love, enjoy strawberries and cream, make someone fall in love with it, learn from experience, use words properly, be the subject of its own thought, have as much diversity of behavior as man, do something really new.*

Some of these are rather **easy** to be replicated by AI. However, some **are not possible**

Overall, programs exceed human performance in some tasks and lag behind on others.

The one thing that it is clear they can't do is *be exactly human*.

# The Limits of AI

The mathematical objection

Turing (1936) and Gödel (1931) proved that certain mathematical questions are in principle **unanswerable** by *particular formal systems*.

Gödel sentence $G(F)$ with the following properties:
- $G(F)$ is a sentence of $F$, but cannot be proved within $F$.
- If $F$ is consistent, then $G(F)$ is true.

Philosophers such as J. R. Lucas (1961) have claimed that this theorem shows that machines are **mentally inferior to humans**:
- machines are **formal systems** that are limited by the incompleteness theorem
- cannot establish the truth of their own Gödel sentence
- Problems with Lucas' claim:
  - Example sentence which **cannot consistently assert** by human else contradiction: *"Lucas cannot consistently assert that this sentence is true."*
  - **No entity**—human or machine—can prove things that are impossible to prove
  - **incompleteness theorem** technically applies only to **formal systems** that are powerful enough to do arithmetic.

# The Limits of AI

## Measuring AI

- Whether machines can pass a behavioral test, which has come to be called the **Turing test**

- The test requires a program to have a conversation (via typed messages) with an interrogator for *five minutes*

- ELIZA program and Internet chatbots such as MGONZ and NATACHATA

- *Eugene Goostman* (chatbot, 2014) fooled 33% of the untrained amateur judges in a Turing test

- AI researchers who crave competition are more likely to concentrate on playing **chess** or **Go** or **StarCraft II**, or taking **an 8th grade science exam**, or identifying **objects in images**.

# The Limits of AI

Measuring AI

- "On the Measure of Intelligence", François Chollet, 2019.

## Skill acquisition as program synthesis

During it's lifetime, a system will face many situations in which it is required to solve a particular task, it does so by generating a skill program.

The system is evaluated and rewarded when it is able to produce program that achieves an acceptable performance solving the task.

The intelligence of the system is directly connected to its ability to generate sufficiently good task-specific skill programs for the given tasks.

An intelligent system synthesizes a program to solve a specific task.

Pearson

# Can Machines Really Think?

Some philosophers claim that a machine that **acts intelligently** would not be actually thinking, but would be only a **simulation of thinking**

*Turing* argues the **polite convention** that everyone and machine think.

*John Searle* rejects the polite convention

## The Chinese room
1. A human, who understands only English, inside a room that contains a rule book, written in English, and various stacks of paper
2. Pieces of paper containing indecipherable symbols are slipped under the door to the room
3. The human follows the instructions in the rule book, finding symbols in the stacks, writing symbols on new pieces of paper, rearranging the stacks, and so on
4. passed back to the outside world
5. it is given that the human does not understand Chinese

Computers are in essence doing the same thing, so therefore computers generate no understanding

# Can Machines Really Think?

Recently released OpenAI ChatGPT: [ChatGPT: Optimizing Language Models for Dialogue (openai.com)](openai.com)

# Can Machines Really Think?

Recently released OpenAI ChatGPT: [ChatGPT: Optimizing Language Models for Dialogue (openai.com)](#)

## Certificate of Intelligence

Upon completion of the IQ Test at IQTest.com,

**chatGPT**

has achieved the following scores:

General Intelligence Quotient: 83

Arithmetic: 83, Spatial Skill: 72, Logical: 78, Spelling: 86, Short Term Memory: 82

Rote Utilization: 82, Algebraic: 83, General Knowledge: 78, Visual Apprehension: 78

Geometric: 84, Vocabulary: 83, Intuition: 82, Computational Speed: 105

### THE IQ TEST SCORE BELL CURVE
**WECAPABLE.COM**

**68%**
of people fall within this range

**IQ TEST**

| 55 Mentally inadequate 2.3% | 70 Low intelligence 13.6% | 85 Below average 34.1% | 100 Above average 34.1% | 115 High intelligence 13.6% | 130 Superior intelligence 2.1% | 145 Exceptionally gifted .13% |

# The Ethics of AI

Given that AI is a powerful technology, we have a **moral obligation** to use it well, to promote the positive aspects and avoid or mitigate the negative ones.

**Positive aspects examples**

- AI can save lives through improved medical diagnosis, new medical discoveries, better prediction of extreme weather events

- AI can improve lives, Microsoft's *AI for Humanitarian Action* program applies AI to recovering from natural disaster

- AI applications in crop management and food production help feed the world

# AI and Gas Emissions?



Global greenhouse gas emissions by sector

This is shown for the year 2016 – global greenhouse gas emissions were 49.4 billion tonnes CO₂eq.

Our World in Data

# The Ethics of AI

**Negative aspects example:**
**Lethal autonomous weapons**

The UN defines a lethal autonomous weapon as one that locates, selects, and engages (i.e., kills) *human targets without human supervision.*

Israel's **Harop** missile is a "loitering munition" with a ten-foot wingspan and a fifty-pound warhead. It searches for up to six hours in a given geographical region for any target that meets a given criterion and then destroys it.

# The Ethics of AI

**Negative aspects example:**
**Lethal autonomous weapons**

Autonomous weapons have been called the *"third revolution in warfare"* after gunpowder and nuclear weapons. Their military potential is obvious...

The debate over autonomous weapons includes **legal, ethical and practical aspects.**
**Legal:** requires the possibility of discriminating between combatants and non-combatants, the judgment of military necessity for an attack, and the assessment of proportionality between the military value of a target and the possibility of collateral damage.



Origin: Israel    Unveiled: 2009

**IAI HAROP**

Producer:
Israel Aerospace Industries

LOITERING MUNITION

Length: 2.5m
Wingspan: 3m

Weight: 135kg
Warhead weight: 25kg

Range: 1,000km
Endurance: 6hours

Features:
- Low-observability
- can autonomously home in on radio emissions
- can be operated in a human-in-the-loop mode

SOOUTHFRONT.ORG

# The Ethics of AI

**Negative aspects example:** **Lethal autonomous weapons**

**Ethical**: some find it simply **morally unacceptable** to delegate the decision to kill humans to a machine.

More than 140 NGOs in over 60 countries are part of the Campaign to **Stop Killer Robots**, and an open letter organized in 2015 by the *Future of Life Institute* organized an <u>open</u> letter was signed by over 4,000 AI researchers and 22,000 others.

**Reliability:** a very serious concern for military commanders, who know well the complexity of battlefield situations. Cyberattacks against autonomous weapons could result in *friendly-fire casualties*

**Practical**: the scale of an attack that can be launched is proportional to the amount of hardware one can afford to deploy.

AI is a **dual use technology**: AI technologies that have peaceful applications can easily be applied to military purposes

# The Ethics of AI



MICRO DRONES KILLER ARMS ROBOTS - AUTONOMOUS ARTIFICIAL
INTELLIGENCE - WARNING !! - YouTube

# The Ethics of AI

Surveillance, security, and privacy

- As of 2018, there were as many as **350 million surveillance cameras** in China and **70 million** in the United States.

- As more of our institutions operate online, more vulnerable to **cybercrime** and **cyberterrorism**. Attackers can use automation to probe for insecurities and they can apply reinforcement learning for phishing attempts and automated blackmail

- Defenders can use unsupervised learning to detect anomalous incoming traffic patterns and various machine learning techniques to detect fraud

- More data on us is being collected by governments and corporation

# The Ethics of AI

Surveillance, security, and privacy

- In the U.S., the *Health Insurance Portability and Accountability Act* (HIPAA) and the *Family Educational Rights and Privacy Act* (FERPA) protect the privacy of medical and student record

- **De-identification**: eliminating personally identifying information (such as name and social security number) so that medical researchers can use the data to advance the common good
  - Federated learning

- **Secure aggregation:** central server doesn't need to know the exact parameter value from each distributed user

# The Ethics of AI

## Fairness and bias

- Machine learning models (especially) can perpetuate societal bias

- Designers of machine learning systems have a **moral responsibility** to ensure that their systems are fair

- Six of the most commonly-used concepts for fairness:
  - Individual fairness
  - Group fairness
  - Fairness through unawareness
  - Equal outcome
  - Equal opportunity
  - Equal Impact

- **COMPAS** is a commercial system for recidivism (re-offense) scoring. It assigns to a defendant in a criminal case a risk score, which is then used by a judge to help make decisions
  - *does not achieve equal opportunity*: the proportion of those who did not re-offend but were falsely rated as high-risk was 45% for blacks and 23% for whites

# The Ethics of AI

Fairness and bias

- **sample size disparity** can lead to biased results.

- In most data sets there will be fewer training examples of minority class

- Machine learning algorithms give better accuracy with more training data, so that means that members of minority classes will experience lower accuracy

- A constrained model may not be able to simultaneously fit both the majority and minority class

- Bias can also come into play in the software development process

- De-bias the data: over-sample from minority classes to defend against sample size disparity

# The Ethics of AI

**Set of best practices**

- Make sure that the software engineers **talk with social scientists** and **domain experts** to understand the issues and perspectives, and consider fairness from the start.

- Create an environment that fosters the development of a **diverse pool of software engineers** that are representative of society.

- **Define what groups your system will support**: different language speakers, different age groups, different abilities with sight and hearing, etc.

# The Ethics of AI

**Set of best practices**

- Optimize for an objective function that **incorporates fairness**.

- **Examine your data** for prejudice and for correlations between protected attributes and other attributes.

- Understand **how any human annotation of data is done**, design goals for annotation accuracy, and verify that the goals are met.

- Don't just track overall metrics for your system; make sure you **track metrics for subgroups** that might be victims of bias.

- **Include system tests that reflect the experience of minority group users**.

- **Have a feedback loop** so that when fairness problems come up, they are dealt with

# The Ethics of AI

**Trust and transparency**

- People need to be able to *trust the systems they use*

- Engineered systems must go through a **verification** and **validation** (V&V) process
  - Verification means that the product satisfies the specifications
  - Validation means ensuring that the specifications actually **meet the needs** of the user and other affected parties

- Certification and safe standards, ISO in other industries

- The AI industry is not yet at this level of clarity, although there are some frameworks in progress, such as IEEE P7001, a standard defining ***ethical design for artificial intelligence and autonomous systems***

- **Transparency**: consumers want to know what is going on inside a system, and that the system is not working against them, whether due to intentional malice, an unintentional bug, or pervasive societal bias that is recapitulated by the system

# The Ethics of AI

Trust and transparency

An AI system that can explain itself is called **explainable AI** (XAI).

A good explanation has several properties:
- it should be understandable and convincing to the user
- it should accurately reflect the reasoning of the system
- it should be complete,
- it should be specific in that different users with different conditions or different outcomes should get different explanations.

The future of work

- An **immediate reduction in employment** when an employer finds a mechanical method to perform work previously done by a person

- **More automation with physical robots**, first in controlled warehouse environments, then in more uncertain environments, building to a significant portion of the marketplace by around 2030.

- The ratio between workers and retirees changes. In 2015 there were less than 30 retirees per 100 workers; by 2050 there may be over 60 per 100 workers

- problems due to **the pace of change**

# The Ethics of AI

<span style="color:purple">Robot Rights</span>

- If robots can feel pain, if they can dread death, if they are considered "persons," then the argument can be made that **they have rights and deserve to have their rights recognize**

- If robots have rights, then they should not be enslaved, and there is a question of whether reprogramming them would be a kind of enslavement

- Another ethical issue involves **voting rights**: a rich person could buy thousands of robots and program them to cast thousands of votes—should those votes count?

- Ernie Davis argues for avoiding the dilemmas of robot consciousness **by never building robots that could possibly be considered consciousness**

# The Ethics of AI

AI Safety

- Design a robot to have low impact, instead of just maximizing utility, maximize the utility minus a weighted summary of all changes to the state of the world.

- Victoria Krakovna (2018) has cataloged **examples of AI agents that have gamed the system**, figuring out how to maximize utility without actually solving the problem that their designers intended them to solve.

    - *Genetic algorithm operating in a simulated world was supposed to evolve fast-moving creatures but in fact produced creatures that were enormously tall and moved fast by falling over.*

- Designers of agents should be aware of these kinds of **specification failures** and take steps to avoid them.

- Need to be very careful in specifying what we want, because with *utility maximizers we get what we actually asked for*. The **value alignment problem**

# The Ethics of AI

EU Regulations

- Trustworthy AI Guidelines, 2019

- [EU AI Act - EU Artificial Intelligence Act](#), 2023

  - Safeguards agreed on *general purpose artificial intelligence*
  - Limitation for the of use biometric identification systems by law enforcement
  - Bans on social scoring and AI used to manipulate or exploit user vulnerabilities
  - Right of consumers to launch complaints and receive meaningful explanations
  - Fines ranging from 35 million euro or 7% of global turnover to 7.5 million or 1.5% of turnover

Hype Cycle for Emerging Technologies, 2024

# AI Components

Sensors and actuators

- Direct access to the world has been absent

- AI systems were built in such a way that humans had to supply the inputs and interpret the outputs.

- The demand for better image processing in cellphone cameras has given us inexpensive high-resolution cameras for use in robotics

- MEMS (micro-electromechanical systems) technology has supplied miniaturized accelerometers, gyroscopes, and actuators small enough to fit in artificial flying insects

- AI systems are at the cusp of moving from ***primarily software-only systems to useful embedded robotic system***

- **Intelligent robots** will first make strides in industry before the home market

# AI Components

Representing the state of the world

- Keeping track of the world requires perception as well as **updating of internal representations**

- Current filtering and perception algorithms can be combined to do a reasonable job of recognizing objects and reporting low-level predicates

- Future progress will require techniques that **generalize to novel situations** without requiring exhaustive examples

- Word embeddings and similar representations can free us from the strict bounds of concepts defined by necessary and sufficient condition

- *However, it remains a **daunting task** to define **general, reusable representation schemes** for complex domains.*

# AI Components

### Selecting Actions

- The primary difficulty in action selection in the real world is coping with **long-term plans**
- Humans apply hierarchical structure on behavior
- Hierarchical reinforcement learning has succeeded in combining these ideas with the MDP formalism
- These methods have not been extended to the partially observable case (POMDPs).

### Deciding What we Want

- The task of picking *the right utility function* is a challenging problem
- We don't have much experience with building complex real-world preference models, let alone probability distributions over such models.
- **Inverse reinforcement learning** is one approach when we have an expert who can perform a task, but not explain it.
- We need better ways of saying what we want and better ways for robots to interpret the information we provide.
- Powerful ecosystem for aggregating user preferences but fail to provide an easy way of opting out
- In the future we will have *personal agents that stick up for our true long-term interests*

# AI Components

Learning

- Current algorithms can cope with quite large problems, reaching or exceeding human capabilities in many task

- Learning can stall when data are sparse, or unsupervised, or when we are dealing with complex representation

- Need advances in transfer learning so that we can take advantage of data in one domain to improve performance on a related domain

- The vast majority of machine learning research today assumes a factored representation

# AI Components

Resources

- Machine learning research and development has been accelerated by the increasing availability of data, storage, processing power, software, trained experts, and the investments needed to support them

- Hundreds of high-quality data sets are available for a range of tasks in computer vision, speech recognition, and natural language processing.

- There is a possibility that **quantum computers** could accelerate AI.

- Currently there are some fast quantum algorithms for the linear algebra operations used in machine learning

- Current quantum computers handle only a few tens of bits, whereas machine learning algorithms often handle inputs with millions of bits and create models with hundreds of millions of parameters.

# AI Architectures

- AI has long had a split between **symbolic systems** (based on logical and probabilistic inference) and **connectionist systems** (based on loss minimization over a large number of uninterpreted parameters).

- **Anytime algorithms**: whose output quality improves gradually over time, so that it has a reasonable decision ready whenever it is interrupted

- **Decision-theoretic meta-reasoning**: applies the theory of information value to the selection of individual computations

- **Meta-reasoning techniques** can be used to design better search algorithms and to guarantee that the algorithms have the anytime property

- **Reflective architecture**: an architecture that enables deliberation about the computational entities and actions occurring within the architecture itself.

# AI Architectures

General AI

- Much of the progress in AI in the 21st century so far has been guided by competition on narrow tasks, such as **the DARPA Grand Challenge for autonomous cars**, the **ImageNet object recognition competition**

- Continued  work on specific tasks (or on individual components) **will not be enough** to reach mastery on a wide variety of tasks

- AI as a field has made a reasonable exploration/exploitation tradeoff, assembling **a portfolio of components**, improving on particular tasks, while also exploring promising and sometimes far-out new ideas.

- Work on components can spur new ideas; for example, generative adversarial networks (GANs) and transformer language models each opened up new areas of research.

# AI Architectures

AI Engineering

- The AI industry **has not yet reached that level of maturity**.

- We do have a variety of powerful tools and frameworks, such as TensorFlow, Keras, PyTorch, CAFFE, Scikit-Learn and SCIPY.

- But many of the most promising approaches, such as GANs and deep reinforcement learning, have proven to be difficult to work with—**they require experience and a degree of fiddling** to get them to train properly in a new domain

- Start with a single huge system and, for each new task, extract from it the parts that are relevant to the task (Jeff Dean)

# AI Architectures

<span style="color:purple">The Future</span>

- AI seems to fit in with other **powerful revolutionary technologies** such as printing, plumbing, air travel, and telephony.

- AI is different from previous revolutionary technologies.

- Improving printing, plumbing, air travel, and telephony to their logical limits would not produce anything to **threaten human supremacy** in the world.

- Improving AI to its current limits certainly could.

In conclusion, AI has made great progress in its short history, but the final sentence of Alan Turing's (1950) essay on *Computing Machinery and Intelligence* is still valid today:

*We can see only a short distance ahead,*
*but we can see that much remains to be done.*

## Chapter 33
## Future Progress in Artificial Intelligence:
## A Survey of Expert Opinion

**Vincent C. Müller and Nick Bostrom**

**Abstract** There is, in some quarters, concern about high–level machine intelligence and superintelligent AI coming up in a few decades, bringing with it significant risks for humanity. In other quarters, these issues are ignored or considered science fiction. We wanted to clarify what the distribution of opinions actually is, what probability the best experts currently assign to high–level machine intelligence coming up within a particular time–frame, which risks they see with that development, and how fast they see these developing. We thus designed a brief questionnaire and distributed it to four groups of experts in 2012/2013. The median estimate of respondents was for a one in two chance that high-level machine intelligence will be developed around 2040–2050, rising to a nine in ten chance by 2075. Experts expect that systems will move on to superintelligence in less than 30 years thereafter. They estimate the chance is about one in three that this development turns out to be 'bad' or 'extremely bad' for humanity.

[Preview of "10.1007-978-3-319-26485-1_33.pdf" (inaoep.mx)](inaoep.mx)

# Summary

- **Philosophers** use the term **weak AI** for the hypothesis that machines could possibly behave intelligently, and **strong AI** for the hypothesis that such machines would count as having actual minds (as opposed to simulated minds)

- AI is a powerful technology, and as such it **poses potential dangers**, through lethal autonomous weapons, security and privacy breaches, unintended side effects, unintentional errors, and malignant misuse. Those who work with AI technology have an ***ethical imperative to responsibly reduce those dangers***.

- AI systems must be able to demonstrate they are **fair**, **trustworthy**, and **transparent**

- There are **multiple aspects of fairness**, and it is impossible to maximize all of them at once. So, a first step is to decide what counts as fair

- **Automation** is already changing the way people work. **As a society, we will have to deal with these changes**.

# What's Next?

1. Projects Presentations: **10-12**, **11-12**
   - One slides deck for group

2. Oral Exam:
   - Only students that have already submitted the report (eng) **at least 1 week before the oral exam date**

3. After the exam I can support you in evolving your project into a **master dissertation** / supervise you on a different project
   - [Thesis Topics (Master Degree) (unipi.it)](#)

4. Interested in a Period Abroad, Internships or PhD application?
   - E-mail me.