

# PHILOSOPHY OF ARTIFICIAL INTELLIGENCE

Fall 2022 Syllabus

---

*Instructor:* Ignacio Ojea Quintana.

*Meeting Time:* Tuesdays 12am to 2pm

*Location:* [Geschw.-Scholl-Pl. 1 \(F\) - F 007 \(Geschossplan\)](#)

*Office Hours:* TBU

*Contact:* TBU

*Uni Link:* [Link](#)

## *Course Description*

---

In this course we will examine central philosophical issues around Artificial Intelligence. The course will have two parts.

The first part will focus on artificial intelligence and the mind. Some of the guiding questions are:

- How can we establish that a computer is intelligent?
- More broadly, how would a theory of the mental in computational terms would look like?
- What problems do computational theories of the mind have?
- What is the difference between symbolic AI and contemporary approaches?
- What can we say about consciousness and artificial intelligence?

In order to answer these questions we will study the Turing Test, Turing Machines, Computationalism and its challenges, Connectionism, and some issues around consciousness, super-intelligence and the Computer Simulation Hypothesis.

The second part of the course will focus on moral and social issues around AI. Some of the guiding questions are:

- What kind of risks does AI pose for humanity?
- What is the moral status of AI?
- How does AI affect fairness relations in society? Can AI be fair?
- AI and Machine Learning applications can be opaque, in the sense that when thing go wrong it is hard to say what failed. How can we improve on AI transparency?
- What role does AI have in popular culture?

We will examine all of these questions by looking at journal articles and chapter selections

## *Prerequisites*

---

The course will have both a technical component and a philosophy component. Mathematics without philosophy is empty, philosophy without mathematics is blind. Students are expected to understand the basics of logic as well as having the basic ability to write a philosophical essay. Having taken an introductory logic course *and* an introductory philosophy course should suffice.

## *Required Texts*

---

All texts will be made available digitally.  
But you should definitely own a printed version of Descartes *Meditations on First Philosophy*.

## *Grading*

---

The course will have two evaluation instances near the end of the semester.

First, a quiz/exam about the content of the class. A list of questions will be provided in sufficient advance, and the exam will consist in a subset of those questions. It is worth 30% of the total grade.

Second, an essay worth 70% of the grade. A guideline with research questions and notes on how to write an essay will be provided.

## *COURSE POLICIES*

---

TBU

## *TENTATIVE COURSE SCHEDULE*

---

Please note that all readings and due dates are subject to change.

Please do the readings **before attending to class**.

**Note: Dates will be corrected according to the academic calendar. This course is about twelve weeks long, factoring in public holidays.**

### **Part 1: Artificial Minds**

**Week 1:** *The Turing Test*

**Tu 18/10**

*Required:*

- A. Turing, "Computing Machinery and Intelligence". *Mind*, Volume LIX, Issue 236, October 1950, Pages 433–460.

*Optional:*

- J.H. Moor, "An Analysis of the Turing Test", *Philosophical Studies*, Vol. 30, No. 4 1976, pp. 249-257
- M. Halina, "Insightful artificial intelligence", *Mind & Language* 36(12), 2021.

**Week 2:** *Turing Machines: Universal Turing Machines, Halting Problem, Church-Turing Thesis.* **Tu 25/10**

*Required:*

- G.S. Boolos, J.P. Burgess, R.C. Jeffrey - *Computability and Logic*, 2002, Cambridge University Press - **Chapter 3**.
- Check this out: <https://turingmachine.io/>

*Optional:*

- G.S. Boolos, J.P. Burgess, R.C. Jeffrey - *Computability and Logic*, 2002, Cambridge University Press - **Chapter 4.**
- G.S. Boolos, J.P. Burgess, R.C. Jeffrey - *Computability and Logic*, 2002, Cambridge University Press - **Chapter 8.**

**Week 3:** *Computationalism*

**Tu 1/11**

*Required:*

- H. Putnam, “The Nature of Mental States”.

*Optional:*

- M. Rescorla, *The Computational Theory of Mind*, Stanford Encyclopedia of Philosophy. ([link](#))
- J. Fodor, “Propositional Attitudes”, *The Monist* 61 (October):501-23 (1978).

**Week 4:** *Challenges to Computationalism*

**Tu 8/11**

*Required:*

- J. Searle, “Can Computers Think?”, In David J. Chalmers (ed.), *Philosophy of Mind: Classical and Contemporary Readings*. Oup Usa (2002).

*Optional:*

- M. Boden, “Escaping from the Chinese Room”, In John Heil (ed.), *Computer Models of Mind*. Cambridge University Press (1988)
- D. Dennett, “Can Machines Think?”, In book: Alan Turing: Life and Legacy of a Great Thinker, 2004.

**Week 5:** *Connectionism and Artificial Neural Networks*

**Tu 15/11**

*Required:*

- P. Smith Churchland, excerpts from Chapter 7 of *Brain-Wise*, ”How do Brains Represent?”, MIT Press, 2002.

*Optional:*

- Notes on Neural Networks (to be uploaded)
- A. Clark, “Connectionism, Competence, and Explanation”, *The British Journal for the Philosophy of Science*, Vol. 41, 1990.
- D. E. Rumelhart, ”The architecture of mind: a connectionist approach’, in M. I. Posner (Ed.), *Foundations of cognitive science* (pp. 133–159). The MIT Press, 1989.

**Week 6:** *Consciousness and AI*

**Tu 22/11**

*Required:*

- R. Descartes, *Meditations* 1st & first half of 2nd. Any translation is good, I suggest you get this book if you do not have it.
- D. Chalmers, “Facing up the problem of consciousness.”, opinion article in *Scientific American*, 1995.

*Optional:*

- T. Nagel, "What is it like to be a bat", *Philosophical Review*, Vol. 83, 1974.
- F. Jackson, "Epiphenomenal Qualia", *The Philosophical Quarterly*, Vol. 32, 1982.

**Week 7: Superintelligence and The Computer Simulation Hypothesis**

**Tu 29/11**

*Required:*

- D. Chalmers, "The Singularity: A Philosophical Analysis." *Journal of Consciousness Studies*, 17(9-10) Excerpt, pp. 1–15 & 19–56, 2010.
- N. Bostrom, "Are We Living in a Computer Simulation?" *Philosophical Quarterly*, 53 (211), pp. 243–255, 2003.

*Optional:*

- J. Prinz. "Singularity and Inevitable Doom", *Journal of Consciousness Studies*, 19 (7-8):77–86, 2012

**Part 2: Moral and Social Issues Around AI**

**Week 8: Existential Risk**

**Tu 6/12**

*Required:*

- V. Müller & M. Cannon, "Existential risk from AI and orthogonality: Can we have it both ways?", *Ratio*, 2021.

*Optional:*

- N. Bostrom, "The Superintelligent Will: Motivation and Instrumental Rationality in Advanced Artificial Agents", *Minds and Machines*, 22(2), pp. 71–85, 2012.
- S. Armstrong, "General Purpose Intelligence: Arguing the Orthogonality Thesis", *Analysis and Metaphysics*, 12, pp. 68–84, 2012.

**Week 9: The Moral Status of AI**

**Tu 13/12**

*Required:*

- D.G. Johnson (2006), "Computer systems: Moral Entities But Not Moral Agents", *Ethics and Information Technology* 8(4), pp. 195–204, 2006.

*Optional:*

- J.P. Sullins, "When is a Robot a Moral Agent?", *International Review of Information Ethics*, 6 (12), pp.23-30, 2006.
- M.A. Warren. *Moral Status: Obligations to Persons and Other Living Things*, Clarendon Press, Excerpt pp. 4–17, 1997.

**Week 10: AI Fairness**

**Tu 20/12**

*Required:*

- B. Hedden, “On statistical criteria of algorithmic fairness”, *Philosophy and Public Affairs*, 49 (2):209-231, 2021.

*Optional:*

- TBU

**Week 11:** *Transparency and Data Bias*

**Tu 10/1**

*Required:*

- M. Günther & A. Kasirzadeh, “Algorithmic and human decision making: for a double standard of transparency”, *AI and Society*, 37 (1):375-381, 2022.

*Optional:*

- Zerilli et al. “Transparency in algorithmic and human decision-making: Is there a double standard?” *Philosophy & Technology* (2019) 32.

**Week 12:** *Artificial Intelligence in Pop Culture* (Terminator, 2001 Space Odyssey, *Her*, *Ex Machina*. Also *Asimov*, *The Matrix*). [For now we use Asimov but we might change that] **Tu 17/1**

*Required:*

- Isaac Asimov, *The Bicentennial Man*. In *The Bicentennial Man and Other Stories*, Doubleday, pp. 138–172, 1976.
- S.L. Anderson, ”The Unacceptability of Asimov’s Three Laws of Robotics as a Basis for Machine Ethics”, In M. Anderson & S. Anderson (eds), *Machine Ethics*, Cambridge University Press, pp. 285, 2011.

**Examination Week:**

**Final Quiz Due**

**Final Paper Due**