

**PHIL GR9180: Topics in Moral Philosophy**  
**Approaches to Applied Ethics—Philosophy of AI**  
**Fall 2024, Monday 2:10-4pm**  
**Philosophy Hall 716**

**1. Course Description**

The philosophy of AI is an emerging field. Right now, AI is importantly concerned with LLMs. It is also concerned with the relation between natural and artificial intelligence. Researchers and public discourse ask whether AI can be “aligned” with values. Accordingly, key questions in Philosophy of AI relate to language, thought, and values.

The seminar starts with the widely debated alignment problem (Part I: weeks 1-4). Independent of how alignment works, it is by no means clear what the desired outcome is. People disagree about values. With which values should AI be aligned? At times the answer is: with “human values.” Are “human values,” in this context, the different and incompatible sets of values human beings have? If yes, what about the values that we should have?

AI researchers often focus on fairness, typically understood as the elimination of bias (Part II: weeks 5-6). We examine relevant notions of fairness and ask how fairness relates to other values, including and especially accuracy.

Next, we ask whether it makes sense to ascribe beliefs and intentions to AIs (Part III: weeks 7-11). Can AIs engage in reasoning, lie and be held responsible? We discuss “explainable AI,” asking whether AI-outputs can be understood.

Finally, we examine questions about language as they apply to AIs (Part IV: weeks 12-14). How can LLMs cope with famously tricky components of language and thought, such as generics and implicature?

Part of the seminar are workshop sessions associated with the ValuesLab. Invited guest speakers come from a range of fields. This seminar aims to contribute to dialogue between philosophers and AI developers, and to a shared vocabulary.

**2. Outline of Readings and Topics**

**Part I: Introduction**

*As an introduction to the philosophy and ethics of AI, we cover questions such as “what is value alignment?” and “how can/do natural and artificial intelligences learn values?”*

**Week 1, Sept 9: Introduction Philosophy of AI and Ethics of AI**

Readings: Peter Railton, “Ethical Learning: Natural and Artificial,” in Matthew Liao (2020); Raphaël Millière, “The Alignment Problem in Context,” arXiv (2023).

Two introductory videos and a historical sketch:

<[https://www.youtube.com/watch?v=zjkBMFhNj\\_g](https://www.youtube.com/watch?v=zjkBMFhNj_g)>

<<https://www.youtube.com/watch?v=ISkAkiAkK7A>>

Bruce Buchanan, [A \(Very\) Brief History of Artificial Intelligence](#)

## **Part II: Alignment and Value Pluralism**

*Debates about “alignment” of AI with “human values” try to recognize that there is no single set of values that people endorse. How should we respond to this? Should AI applications be tailored to the views users already have? We work toward an updated “map” of outlooks in normative ethics, as a starting point for questions about alignment, a plurality of evaluative outlooks, and questions about relativism and universal values.*

### **Week 2, Sept 16: Familiar Options in Normative Ethics—Teleology versus Deontology, Consequentialism versus Kantian Ethics**

Readings: John Rawls on teleology and deontology, “Classical Utilitarianism,” Part I ch. 5, *A Theory of Justice* (1971); Jens Timmermann, “What’s Wrong with ‘Deontology’?” *Proceedings of the Aristotelian Society*, Volume 115 (2015); Mill, *Utilitarianism II*; Shelly Kagan, *The Limits of Morality* (OUP 1989), selection.

### **Week 3, Sept 23: Recent Additions to the “Map” of Options in Normative Ethics—Virtue Ethics, Confucius, Buddhist ethics, Cognitive Science, Philosophy of Race, Ethics of Care**

Readings: Julia Annas, “Ancient Eudaimonism and Modern Morality,” in: Bobonich *Cambridge Companion to Ancient Ethics* (CUP, 2017); Katja Maria Vogt, *Desiring the Good* (OUP 2017), chs 3 and 5; Amber Carpenter and Pierre-Julien Harter, “Introduction” in *Crossing the Stream, Leaving the Cave* (OUP 2024); SEP “Implicit Bias” <<https://plato.stanford.edu/entries/implicit-bias/>>; Virginia Held, *The Ethics of Care: Personal, Political, and Global*. Oxford University Press, 2005, pp. 9-28; Carol Gilligan, *In a Different Voice: Psychological Theory and Women’s Development*, Harvard, 1982. 1-11; Charles Mills “White Ignorance” In Shannon Sullivan & Nancy Tuana eds., *Race and Epistemologies of Ignorance*, State Univ of New York (2007), 11-38. Cf. Mills, “Global White Ignorance,” in: M. Gross & L. McGoey (Eds.), *Routledge international handbook of ignorance studies*, Routledge/Taylor & Francis Group (2015), 217-227.

### **Week 4, Sept 30: Alignment and Decisions**

Readings: Ruth Chang, “Humans in the Loop,” in ed. David Edmonds, *AI Morality* (OUP 2024/ forthcoming); Edna Ullmann Margalit, “Big Decisions, Opting, Converting, Drifting” *Royal Institute of Philosophy Supplements* 58 (2006): 157-172; Oliver Klingefjord, Ryan Lowe, Joe Edelman, “What are human values, and how do we align AI to them?” (2024) <<https://arxiv.org/abs/2404.10636>>. Background readings on pluralism: Jeremy Waldron, “Public Reason and “Justification” in the Courtroom,” *Journal of Law, Philosophy and Culture* (2007): 107-134. Elinor Mason, “Value Pluralism.” *The Stanford Encyclopedia of Philosophy* (Summer 2023 Edition), Edward N. Zalta & Uri Nodelman (eds.), URL = <<https://plato.stanford.edu/archives/sum2023/entries/value-pluralism/>>.

## **Part III: Fairness and Accuracy**

*AI researchers seem to think of fairness as the opposite of bias. We discuss how fairness relates to accuracy, justice, and other values.*

### **Week 5, Oct 7: Fairness**

Readings: Richard Zemel et al, “Learning Fair Representations” (2013); Kate Vredenburg, “Fairness,” in Justin Bullock (ed.), *The Oxford Handbook of AI Governance*, OUP 2022; Drago Plečko and Elias Bareinboim (2024), “Causal Fairness Analysis”, *Foundations and Trends in Machine Learning*: Vol. 17, No. 3, pp 1–238. DOI: 10.1561/2200000106. Additional reading: Talia Gillis, Bryce McLaughlin, Jann Spiess “On the Fairness of Machine-Assisted Human Decisions” (2023) [this paper looks at uses of AI models where human decision-makers work with input from an AI]

### **Part III: Mental States?**

*It is tempting to speak of what an AI “believes,” “intends,” and so on. Whether these ascriptions are plausible depends, in part, on what we think beliefs, intentions, and so on, are. People often speak of AIs as a black box, signaling that it seems opaque how an AI arrives at its output. We ask what explainability amounts to, and how it relates to traditional ideas in philosophy of mind and epistemology.*

### **Week 6, Oct 14: Belief-Ascriptions?**

Readings: Eric Schwitzgebel, “Belief,” *The Stanford Encyclopedia of Philosophy* (Spring 2024 Edition), Edward N. Zalta & Uri Nodelman (eds.), URL = <<https://plato.stanford.edu/archives/spr2024/entries/belief/>>, Section 1; Bernard Williams, “Deciding to Believe” (1973); Murray Shanahan, “Talking about Large Language Models,” (2023); Harvey Lederman and Kyle Mahowald, “Are Language Models More Like Libraries or Like Librarians? Bibliotechnism, the Novel Reference Problem, and the Attitudes of LLMs” (2024).  
Additional readings: Timothy Williamson, “Is Knowing a Mental State?” (1995) and “Knowledge First Epistemology,” (2011).

### **Week 7, Oct 21, AI Agents**

Readings: Katja Maria Vogt, “The Guise of the Good,” *Desiring the Good* (OUP 2017), chapter 5; Beba Cibralic and James Mattingly, “Machine agency and representation,” (2021); Chad Firestone, “Performance vs. competence in human-machine comparisons,” *PNAS* (2020).  
Background: Kieran Setiya, “Intention,” *The Stanford Encyclopedia of Philosophy* (Winter 2022 Edition), Edward N. Zalta & Uri Nodelman (eds.), URL = <<https://plato.stanford.edu/archives/win2022/entries/intention/>>; Pierre Jacob, “Intentionality,” *The Stanford Encyclopedia of Philosophy* (Spring 2023 Edition), Edward N. Zalta & Uri Nodelman (eds.), URL = <<https://plato.stanford.edu/archives/spr2023/entries/intentionality/>>

### **Week 8, Oct 28: Lying and Hallucinating**

Readings: James Mahon, “The Definition of Lying and Deception,” *SEP* (2015) <<https://plato.stanford.edu/entries/lying-definition/>>; Andreas Stokke, “Lying and Asserting,” *Journal of Philosophy* (2013); Karen Lewis, “Language: Gricean Pragmatics” <<https://www.youtube.com/watch?v=we6uSVf4qss>>; Levinstein, B. A., & Herrmann, D. A. (Accepted/In press). “Still no lie detector for language models: probing empirical and conceptual roadblocks,” *Philosophical Studies*; Kristina Suchotzki and Matthias Gamer, “Detecting deception with artificial intelligence: promises and perils,” *Science & Society* 28.6 (2024), 481-483.

### **Week 9, Nov 11: ValuesLab Workshop on Fairness in AI**

Guest: Joshi Shalmali (CU Medical School) on her work on fairness in AI.

Commentator: TBD

**Week 10, Nov 18: ValuesLab Workshop**

Guest: Phillip Koralus, Oxford, topic TBD

Commentator: TBD

**Week 11, Nov 25: Explainable AI (XAI), Competence, and Reliability**

Readings: The General Data Protection Regulation, issued by the European Union in 2018, speaks of a “right to explanation”; “Is explanation ex post rationalization?” <<https://ykulbashian.medium.com/pragmatics-precedes-semantics-c8fb67a19ead>>; Federica Russo, Eric Schliesser, Jean Wagemans, “Connecting Ethics and Epistemology of AI” (2023); Alan Goldman, “What is Justified Belief?” (selection, 1979); Ernest Sosa, *A Virtue Epistemology* (2007), “Lecture 2: A Virtue Epistemology”; Herman Cappelen and Joshua Dever, *Making AI Intelligible: Philosophical Foundations* (OUP 2021), “Introduction.”

**Part IV: Language**

*Much of the current debate about AI is about LLMs. This should put questions about language front and center in the discussion of ethical AI. We turn to two pervasive dimensions of language—and thought—that have recently received much attention: the semantics of generics and conversational implicature. Similarly, we look at the notion of linguistic communities and its role in philosophy of language and computational linguistics.*

**Week 12, Dec 2: ValuesLab Workshop**

Guest: TBA

Commentator: TBD

**Week 13, Dec 9: Generics and Implicature**

Readings: Emily Allaway et al on generics as incomplete information; “Social and Political Aspects of Generic Language and Speech.” McKeever and Sterken; Karen Lewis, “Gricean Pragmatics,” Laura Ruis, Akbir Khan, Stella Biderman, Sara Hooker, Tim Rocktäschel, Edward Grefenstette, “The Goldilocks of Pragmatic Understanding: Fine-Tuning Strategy Matters for Implicature Resolution by LLMs,” ArXiv (2022).

Background reading: Raphaël Millièrè & Cameron Buckner, “A Philosophical Introduction to Language Models – Part I: Continuity With Classic Debates,” arXiv (2024).

**3. Requirements**

Readings: For each class, we have a range of readings. Students are expected to come to class having read the material carefully and ready to participate in discussion.

Presentations: Students are expected to participate in team presentations (5 minutes per team/paper). Presentations go along with a 1-page handout that clearly lays out key concepts and arguments. Students are expected to run the handout by the instructor at the latest 2 days prior to the presentation. They are invited to come to office hours to discuss the material. *Please get in touch early in the semester if you are interested in contributing to a team-presentation!*

E-credit writing requirements: Regular contributions to in-class presentations and team work in preparing presentations. Two 6-page papers or one 12-page paper. Deadlines for the 2-paper option: Dec 25 and Dec 13. Deadline for the 1-paper option Dec 13.

R-credit requirements: Careful reading in preparation for class and participation in class.

#### **4. Academic Integrity and Honor Code**

Please consult Columbia University's policies on academic integrity as well as Columbia's honor code:

<http://www.college.columbia.edu/academics/academicintegrity>

<https://www.college.columbia.edu/ccschonorcode>

<http://bulletin.columbia.edu/general-studies/undergraduates/academic-policies/academic-integrity-community-standards/>

These policies explain Columbia University's academic regulations and how you can safeguard the integrity of your original work. Plagiarism and other forms of academic dishonesty are serious offenses. Please take the time to familiarize yourself with the details of what constitutes plagiarism and academic dishonesty. You are expected to confirm to these policies in your academic work. It is important that you understand that academic dishonesty can lead to disciplinary action, including failure in the course and suspension, or even expulsion, from the University.

#### **5. Accommodations for Students with Disabilities**

In order to receive disability-related academic accommodations, students must first be registered with Disability Services. More information on the Disability Services registration process is available online at <[www.health.columbia.edu/ods](http://www.health.columbia.edu/ods)>. Registered students must present an Accommodation Letter to the professor before an exam or other accommodations can be provided. Students who have, or think they may have, a disability are invited to contact Disability Services for a confidential discussion at (212) 854-2388 (Voice/TTY) or by email at <[disability@columbia.edu](mailto:disability@columbia.edu)>.