

**INLS 690**  
**PHILOSOPHY AND ETHICS OF AI**

Instructor: Fox  
Email: foxm@unc.edu  
Room: Manning 303  
Schedule: M/W 1:25-2:40pm  
Office: Zoom, Greenlaw 505  
Office Hours: 3:30-5:30pm on Wednesdays or by appointment

**DESCRIPTION AND OBJECTIVES**

This course covers such topics as the Turing test and the frame problem in the philosophy of AI, superintelligence and the so-called singularity—when AI becomes uncontrollable and independent, when it governs and improves itself—the existential risk of Artificial General Intelligence (AGI), the moral status of AI, the design of ethical AI, and the possible implications of merging AI with humans. We will engage with leading philosophers and ethicists in the field, reflect on stories by Isaac Asimov and Harlan Ellison, and critique the films *2001: A Space Odyssey* and *Her*. We will also explore what some consider to be myths about AI. Overall, this course will teach you how to think critically about AI from philosophical, ethical, and imaginative perspectives.

**TEXTS**

All readings for this course are available under Course Reserves on Canvas.

**REQUIREMENTS (AND GRADE DISTRIBUTION)**

**Participate (20%)**

This course will be conducted as a seminar, so your participation in class discussion will be integral to its success. Please come to class well-prepared to discuss the readings. Before each class period, you'll submit a brief response to that class period's readings to a Canvas discussion board. You can address part of a reading or the reading as a whole. The purpose of these responses is to help you grasp what you've read and to help prepare the class for in-class discussion.

Your participation grade will be based mostly on the quality of your commentary on Canvas and in-class. Try to have something to say each class period.

I'll provide participation grades at the end of the semester.

**Reflect (30%)**

You'll submit 2 1000-word pieces that will each critically or philosophically engage with one or more of the texts or films included in a particular week's list of materials.

## **Write (40%)**

For your term paper, which should be 2500-3000 words, you have several options:

- 1) An extended critical or philosophical engagement of a topic from the course or a related topic
- 2) A well-researched magazine article in the style of an article from Nautilus Magazine that addresses some aspect of the current or future state of AI and its philosophical, ethical, or cultural implications
- 3) A critical essay on the representation of AI in literature or film that addresses philosophical, ethical, or other cultural concerns

## **Present (10%)**

At the end of the semester, you'll present a lightning talk, about 5 minutes long, on your term paper.

## **GRADING SCALE**

### **Undergraduate**

A	> 92
A-	90-92
B+	88-89
B	83-87
B-	80-82
C+	78-79
C	73-77
C-	70-72
D+	68-69
D	63-67
D-	60-62
F	< 60

### **Graduate**

H	– High Pass
L	– Low Pass
P	– Pass
F	– Fail
I	– Incomplete

## **HONOR CODE**

Your full participation and observance of the University's Honor Code are expected. And you are not permitted to upload any content from this course, including recordings you may be

permitted to make through Accessibility Resources & Service, to the web in any form, including but not limited to Chegg, Course Hero, Coursera, Google Drive, etc. If you post course content, you may be violating my intellectual property rights. If you post your own work from this course, you are allowing sites to profit from your intellectual property. In utilizing web sources to upload or download course content, you risk violating the University's Honor Code.

## ACCESSIBILITY RESOURCES AND SERVICE

The University of North Carolina at Chapel Hill facilitates the implementation of reasonable accommodations, including resources and services, for students with disabilities, chronic medical conditions, a temporary disability or pregnancy complications resulting in barriers to fully accessing University courses, programs and activities. Accommodations are determined through the Office of Accessibility Resources and Service (ARS) for individuals with documented qualifying disabilities in accordance with applicable state and federal laws. See the ARS Website for contact information (<https://ars.unc.edu>) or email [ars@unc.edu](mailto:ars@unc.edu).

## MISCELLANEOUS

I may make changes to the syllabus as necessary over the course of the semester.

## SCHEDULE

A \* indicates a due date.

<b>Jan. 9</b>	<b>Overview and Introductions</b>
<b>Jan. 11</b>	<b>The Turing Test</b> Turing, "Computing Machinery and Intelligence." <i>PEAI</i>
<b>Jan. 16</b>	<b>No class (MLK Day)</b>
<b>Jan. 18</b>	<b>The Chinese Room Thought Experiment 1</b> Searle, "Minds, Brains, and Programs." <i>PEAI</i>
<b>Jan. 23</b>	<b>The Chinese Room Thought Experiment 2</b> Boden, "Escaping from the Chinese Room" <i>PEAI</i>
<b>Jan. 25</b>	<b>Mind as Computer</b> Searle, "Is The Brain a Digital Computer?" <i>Proceedings and Addresses of the American Philosophical Association</i>
<b>Jan. 30</b>	<b>The Frame Problem 1</b> Dennett, "Cognitive Wheels: The Frame Problem of AI." <i>PEAI</i>  Optional reading: Hanahan, "The Frame Problem." <i>The Stanford Encyclopedia of Philosophy</i>
<b>Feb. 1</b>	<b>The Frame Problem 2</b> Searle, "The Background of Meaning, <i>Speech Act Theory and Pragmatics</i>
<b>Feb. 6</b>	<b>Superintelligence and the Singularity 1</b> Chalmers, "The Singularity: A Philosophical Analysis," <i>Journal of Consciousness Studies</i>

<b>Feb. 8</b>	<b>Superintelligence and the Singularity 2</b> Goertzel. "Should Humanity Build a Global AI Nanny to Delay the Singularity Until Its Better Understood?" <i>Journal of Consciousness Studies</i>
<b>Feb. 13</b>	<b>No class (Well-being Day)</b>
<b>*Feb. 15</b>	<b>The Singularity and the Computer Simulation Hypothesis</b> Prinz, "Singularity and Inevitable Doom." <i>Journal of Consciousness Studies</i> <b>*Response Paper 1 due*</b>  Optional reading: Bostrom, "Are We Living in a Computer Simulation?" <i>Philosophical Quarterly</i>
<b>Feb. 20</b>	<b>Artificial General Intelligence (AGI) and Existential Risk 1</b> Bostrom. "The Superintelligent Will: Motivation and Instrumental Rationality in Advanced Artificial Agents. Minds and Machines." <i>Mind and Machines</i>  Optional reading: Omohundro. "The Basic AI Drives." <i>Proceedings of the AGI08 Workshop</i> . Armstrong. "General Purpose Intelligence: Arguing the Orthogonality Thesis." <i>Analysis and Metaphysics</i>
<b>Feb. 22</b>	<b>Artificial General Intelligence (AGI) and Existential Risk 2</b> Häggström. "Challenges to the Omohundro–Bostrom Framework for AI Motivations." <i>Foresight</i>  Optional reading: Goertzel. "Superintelligence: Fears, Promises and Potentials." <i>Journal of Evolution and Technology</i>
<b>Feb. 27</b>	<b>Moral Status of AI 1</b> Asimov. "The Bicentennial Man." <i>The Bicentennial Man and Other Stories</i> .
<b>Mar. 1</b>	<b>Moral Status of AI 2</b> Anderson. "The Unacceptability of Asimov's Three Laws of Robotics as a Basis for Machine Ethics." <i>Machine Ethics</i> .  Optional reading: Michael LaBossiere (2017). Testing the Moral Status of Artificial Beings; or I'm Going to Ask You Some Questions." <i>Robot Ethics 2.0: From Autonomous Cars to Artificial Intelligence</i>
<b>Mar. 6</b>	<b>Ethical AI 1</b> Allen et al. "Prolegomena to Any Future Artificial Moral Agent." <i>Journal of Experimental &amp; Theoretical Artificial Intelligence</i>  Optional reading: Awad et al. "The Moral Machine Experiment." <i>Nature</i>
<b>Mar. 8</b>	<b>Ethical AI 2</b>

	Yampolskiy et al. "Safety Engineering for Artificial General Intelligence." <i>Topoi</i>
<b>Mar. 13</b>	<b>No class (Spring Break)</b>
<b>Mar. 15</b>	<b>No class (Spring Break)</b>
<b>*Mar. 20</b>	<b>Artificial You 1</b> Schneider. Introduction-Chapter 2, <i>Artificial You</i> <b>*Response Paper 2 due*</b>
<b>Mar. 22</b>	<b>Artificial You 2</b> Schneider. Chapters 3-4, <i>Artificial You</i>
<b>Mar. 27</b>	<b>Artificial You 3</b> Schneider. Chapters 5-7, <i>Artificial You</i>
<b>Mar. 29</b>	<b>Artificial You 4</b> Schneider. Chapter 8 and Conclusions, <i>Artificial You</i>
<b>Apr. 3</b>	<b>Workshop</b> Come prepared to discuss ideas for your term paper.
<b>Apr. 5</b>	<b>Literature</b> Ellison, "I Have No Mouth and I Must Scream."
<b>Apr. 10</b>	<b>Film 1</b> <i>2001: A Space Odyssey</i>
<b>Apr. 12</b>	<b>Film 2</b> <i>2001: A Space Odyssey</i>
<b>Apr. 17</b>	<b>Film 3</b> <i>Her</i>
<b>Apr. 19</b>	<b>Film 4 (cont'd)</b> <i>Her</i>
<b>Apr. 24</b>	<b>Presentations 1</b>
<b>Apr. 26</b>	<b>Presentations 2</b>
<b>*May 1</b>	<b>*Term Paper due*</b>