

DRAFT Syllabus: Syllabus will be finalized closer to the start of the teaching semester.

RESPONSIBLE, TRUSTWORTHY, AND SUSTAINABLE ENGINEERING AI (ENAI604) Sections XXXX and XXXX

Term: *Fall 2026*

Professor: Sanghamitra Dutta

Pronouns: she/her

Office Phone: 301-405-2677 Email: sanghamd@umd.edu

Credits: ###

Course Dates: From Aug 31, 2026 - Dec 21, 2026

Course Times: ### Classroom: ###

Canvas/ELMS: ###

Office Hours: ###

Teaching Assistant: ###

Pronouns: ### Email: ###

Office Hours: ###

Course Description

With the growing use of machine learning and generative AI in high-stakes applications, it is of utmost importance that these models are trustworthy, responsible, and sustainable. ENAI604 provides a comprehensive coverage of trustworthy machine learning, equipping students with the knowledge and tools to design and deploy AI systems more reliably and efficiently in practice. The course will deep-dive into four core learning modules:

- (i) Interpretability & Explainability (XAI): Move beyond "black-box" models. Master techniques to clearly articulate Al-driven decision-making for various human stakeholders, enhancing transparency and user confidence.
- (ii) Robustness: Build defenses against an adversarial world. Learn cutting-edge techniques to fortify models against adversarial attacks, data poisoning, jailbreaking, and failure in the face of unseen real-world data.
- (iii) Ethics: Ensure fairness, prevent bias, safeguard data privacy, and master methods to prevent unwanted data leakage. Align your AI systems with core ethical principles and regulatory requirements.
- (iv) Compression & Sustainability: Address the colossal resource demands of modern AI. Explore compression and distillation techniques to deploy large models efficiently in smaller, resource-constrained environments.

Going beyond traditional neural networks, the course will cover **transformers**, large language models (**LLMs**), and vision language models (**VLMs**), with a focus on trustworthy and efficient deployment for real-world applications.

Prerequisites

This course requires prior completion of an introductory AI/ML course (with coverage of matrix operations and basic neural networks) and some experience with Python programming. Advanced AI/ML topics, such as generative AI, will be covered in the course as needed. No prior experience with LLMs or generative AI is required.

Learning Outcomes

After successfully completing this course, you will be able to:

- Apply a broad range of techniques for interpretability, robustness, ethics, and compression in AI systems
- Demonstrate an understanding of advanced machine learning models such as transformers, Large Language Models (LLMs), and Vision Language Models (VLMs) with a focus on trustworthy and efficient deployment for real-world applications
- Analyze and evaluate machine learning models for improved trust and efficiency
- Implement and deliver real-world machine learning models in accordance with specific compliance requirements and/or resource constraints, in a team environment
- Identify potential failure modes in existing machine learning models, propose creative solutions, and analyze their benefits, costs, and limitations

Course Materials

Required Resources

- Book(s):
 - o Interpretable Machine Learning: A Guide for Making Black Box Models

Author: Christoph Molnar

https://christophm.github.io/interpretable-ml-book/

Hands-On Large Language Models: Language Understanding and Generation

Author: Jay Alammar, Maarten Grootendorst

https://www.oreilly.com/library/view/hands-on-large-language/9781098150952/

ISBN-10: 1098150961

ISBN-13: 978-1098150969

- Other online reading materials will be shared as required
- Application/Software. Google Colab (Some assignments may require nominal GPU renting)
- Total estimated costs of required course materials: \$50.00 Textbook, \$50.00 Google Cloud

Supplemental Resources (no purchase required)

- Readings:
 - Build a Large Language Model (From Scratch)

Author: Sebastian Raschka

https://www.manning.com/books/build-a-large-language-model-from-scratch

- Other online reading materials will be shared as required
- Hardware/Software: Regular laptop with Wifi to be able to access Google Colab during class

Course Structure

This course includes both on-campus and online sections. To attend synchronously online, log into ELMS-Canvas at the time of the Section 0101 class [include day/time] and select "Video Conference" from the left side menu. This will open a Zoom link to the live classroom.

For asynchronous online students, all lectures will be recorded and made available on ELMS-Canvas under "Panopto Recordings/Video Lectures" within 24 hours of the class time. Be sure to review the recorded lectures on time.

On-campus students come to class prepared to engage with the lecture and materials. Online students, be sure to log into Canvas regularly and participate in discussions and activities. Regardless of the section you are enrolled in, participation is expected.

Please note that F1 students enrolled in the on-campus section are required to attend in person. If you have a conflict on a particular day, please reach out to me in advance to discuss.

Communication Guidelines

Communicating with the Instructor

My goal is to be readily available to you throughout the semester. I can be reached by email at sanghamd.umd.edu. Please DO NOT email me with questions that are easily found in the syllabus or on ELMS-Canvas (e.g., When is this assignment due? How much is it worth? etc.), but please DO reach out about personal, academic, and intellectual concerns/questions.

While I will do my best to respond to emails within 24 hours, you will more likely receive email responses from me on Fridays from 11 AM to 6 PM EST.

When constructing an email to me please put "ENAI 604 (Section XXXX): Your Topic" in the subject line. This will draw my attention to your email and enable me to respond to you more quickly.

Additionally, please review <u>These tips for 'How to email a professor'</u>. By following these guidelines, you will be ensured to receive a timely and courteous response.

Finally, if you need to discuss issues not appropriate for the classroom and/or an email, we can arrange to talk by phone, over Zoom, or in person. Send me an email asking for a meeting and we can set something up.

Announcements

I will send IMPORTANT messages, announcements, and updates through ELMS-Canvas. To ensure you receive this information in a timely fashion, make sure your email and announcement notifications (including changes in assignments and/or due dates) are enabled in ELMS-Canvas (How to change notification settings in CANVAS).

Log into our ELMs-Canvas course site at least once every 24-hour period to check your inbox and the Announcements page.

Names/Pronouns and Self-Identifications

The University of Maryland recognizes the importance of a diverse student body, and we are committed to fostering inclusive and equitable classroom environments. I invite you, if you wish, to tell us how you want to be referred to in this class, both in terms of your name and your pronouns (he/him, she/her, they/them, etc.). Keep in mind that the pronouns someone uses are not necessarily indicative of their gender identity. Visit trans.umd.edu to learn more.

Additionally, it is your choice whether to disclose how you identify in terms of your gender, race, class, sexuality, religion, and dis/ability, among all aspects of your identity (e.g., should it come up in classroom conversation about our experiences and perspectives) and should be self-identified, not presumed or imposed. I will do my best to address and refer to all students accordingly, and I ask you to do the same for all of your fellow Terps.

Communicating with your Peers

With a diversity of perspectives and experience, we may find ourselves in disagreement and/or debate with one another. As such, it is important that we agree to conduct ourselves in a professional manner and that we work together to foster and preserve a virtual classroom environment in which we can respectfully discuss and deliberate controversial questions. I encourage you to confidently exercise your right to free speech—bearing in mind, of course, that you will be expected to craft and defend arguments that support your position. Keep in mind, that free speech has its limit and this course is NOT the space for hate speech, harassment, and derogatory language. I will make every reasonable attempt to create an atmosphere in which each student feels comfortable voicing their argument without fear of being personally attacked, mocked, demeaned, or devalued.

Any behavior (including harassment, sexual harassment, and racially and/or culturally derogatory language) that threatens this atmosphere will not be tolerated. Please alert me immediately if you feel threatened, dismissed, or silenced at any point during our semester together and/or if your engagement in discussion has been in some way hindered by the learning environment.

Netiquette Policy

Netiquette is the social code of online classes. Students share a responsibility for the course's learning environment. Creating a cohesive online learning community requires learners to support and assist each other. To craft an open and interactive online learning environment, communication has to be conducted in a professional and courteous manner at all times, guided by common sense, collegiality, and basic rules of etiquette.

Grading

Grade Breakdown

Assignment	Percentage %
In-class problem sets and quizzes	10%
4 Homeworks with hands-on data visualization	40%
Midterm Exam	25%
Semester-Long Project	25%
Total	100%

Course Assignments

Homework Assignments

• Four homework assignments will be released throughout the course. Homeworks would require training or fine-tuning various machine learning models, including advanced models like transformers and language models, on real-world datasets and then applying the techniques taught in class to make these models

- responsible, trustworthy, and sustainable. The homework assignments could be solved on Google Colab and might require nominal GPU usage.
- The homework assignments will provide thorough hands-on experience with real-world datasets. Students will gain a crucial appreciation for how machine learning models can be unreliable, and will learn practical techniques to build more responsible, trustworthy, and sustainable models. Students will be able to compare various algorithms and highlight their trade-offs, strengths, and weaknesses.

In-Class Problem Sets & Quizzes

- Every lecture will include 2-3 in-class problem sets. Some problem sets will involve quick data analysis and code completion in Google Colab notebooks shared during class. Some problem sets will involve a small quiz based on the lecture content. Final answers would have to be submitted for grading.
- The in-class problem sets will focus on quick real-time visualization of the covered techniques. Problems will be designed to provide immediate hands-on experience on public datasets for enhanced engagement, retention, and applied understanding beyond theoretical content.

Participation & Engagement

• Students are expected to actively participate in the live lectures by asking questions or volunteering to share their thoughts or course reflections during each lecture. Every student must share their reflections at least once during the course.

Team Project

- The course includes one semester-long project. Students are expected to form teams of 3-5 and implement an existing or new technique on a real-world dataset. The project proposals will be due at the end of November, and the final presentations and report will be due at the end of the course.
- The purpose of the team project is to enable students to deep dive on a sub-topic of interest within the realm of the course in a team environment. Students will critically evaluate their solutions and choose the most suitable tool for different contexts and stakeholders. More information can be found on ELMS.

Midterm Exam

- The exams will be held in class and will include a mix of multiple-choice questions and problems with short answers to be solved on paper.
- The exam is closed-book. Students can bring in a single letter-sized cheat sheet.

Grading of Assignments

All assignments will be graded according to a predetermined set of criteria (i.e., rubric), which will be communicated to students before the assignment is submitted.

To progress satisfactorily in this class, students need to receive timely feedback. To that end, it is my intention to grade all assignments within **10 days** of their due date. If an assignment is taking longer than expected to grade, students will be informed of when they can expect to see their grade.

Grade Computation

All assessment scores will be posted on ELMS/Canvas page. If you would like to review any of your grades (including the exams), or have questions about how something was scored, please email me to schedule a time for us to meet and discuss.

It is expected that you will submit work by the deadline listed in the syllabus and/or on ELMS-Canvas. Late work will be penalized according to the late work policy described in the **Course Policies and Procedures** section below.

Grade Disputes: I am happy to discuss any of your grades with you, and if I have made a mistake, I will immediately correct it. Any formal grade disputes must be submitted in writing and within one week of receiving the grade.

Final letter grades are assigned based on the percentage of total assessment points earned. To be fair to everyone, I have to establish clear standards and apply them consistently, so please understand that being close to a cutoff is not the same as making the cut $(89.99 \neq 90.00)$. It would be unethical to make exceptions for some and not others.

Final Grade Cutoffs

Letter Grade	Cutoff
A+	97%
Α	94%
Α-	90%
B+	87%
В	84%
B-	80%
C+	77%
С	74%
C-	70%
D+	67%
D	64%
D-	60%
F	<60%

Course Schedule

Week #	Торіс	Deliverable
1	Course Overview & Python Bootcamp	
2	Background on Machine Learning Models	In-class problem set
3	Interpretable Machine Learning Part I: From inherently interpretable predictors to posthoc explanations LIME & SHAP	In-class problem set
4	Interpretable Machine Learning Part II: Counterfactuals, Saliency Maps, and Concepts	In-class problem set
5	Introduction to Generative AI Part I: Transformers	In-class problem set, HW 1 due
6	Introduction to Generative AI Part I: LLMs & VLMs	In-class problem set
7	Mechanistic Interpretability for LLMs: An Overview	In-class problem set
8	Midterm Review and In-Class Practice	HW 2 due
9	Midterm	Handwritten paper exam after class
10	Robustness Part I: Traditional Neural Networks	In-class problem set
11	Robustness Part II: Jailbreaking and Uncertainty in LLMs	In-class problem set, Project Proposals due

6 Last updated 11/3/2025

Week #	Торіс	Deliverable
1	Course Overview & Python Bootcamp	
2	Background on Machine Learning Models	In-class problem set
3	Interpretable Machine Learning Part I: From inherently interpretable predictors to posthoc explanations LIME & SHAP	In-class problem set
4	Interpretable Machine Learning Part II: Counterfactuals, Saliency Maps, and Concepts	In-class problem set
5	Introduction to Generative AI Part I: Transformers	In-class problem set, HW 1 due
6	Introduction to Generative AI Part I: LLMs & VLMs	In-class problem set
7	Mechanistic Interpretability for LLMs: An Overview	In-class problem set
8	Midterm Review and In-Class Practice	HW 2 due
12	Ethics Part I: Fairness in Machine Learning	In-class problem set, HW 3 due
13	Ethics Part II: Privacy & Data Leakage	In-class problem set
14	Sustainability: Compression & Distillation	In-class problem set
15	Final Project Presentations	Oral presentation, HW 4 due
16	Final Project Presentations	Oral presentation & Final report

Note: This is a tentative schedule, and subject to change as necessary – monitor ELMS-Canvas for current deadlines. In the unlikely event of a prolonged university closing or an extended absence from the university, adjustments to the course schedule, deadlines, and assignments will be made based on the duration of the closing and the specific dates missed.

Course Policies and Procedures

The University of Maryland's conduct policy indicates that course syllabi should refer to a webpage of course-related policies and procedures. For a complete list of graduate course-related policies, visit the <u>Graduate School website</u>. Below are course-specific policies and procedures that explain how these Graduate School policies will be implemented in this class.

Satisfactory Performance

The Graduate School expects students to take full responsibility for their academic work and academic progress. The student, to progress satisfactorily, must meet all the academic requirements of this course. Additionally, each student is expected to complete all readings and any preparatory work before each class session, come to class prepared to make substantive contributions to the learning experience, and to proactively communicate with the instructor when challenges or issues arise.

Questions about Assignments

Please ask all questions you may have about an assignment by 6:00 PM the day before the assignment is due. Any questions asked after that time may not be answered in time for you to make changes to your work.

Late Work Policy

Assignments should be completed by the due date and time listed with the assignment, on the syllabus, and/or in the course calendar. If you are unable to complete an assignment by the stated due date, it is your responsibility to contact your instructor to discuss an extension, at least 24 hours BEFORE the assignment is due. Extensions are not guaranteed, but may be granted at the instructor's discretion.

Homework assignments will be open for the next 48 hours (2 days). Grade will decrease linearly each hour until it becomes 0. Work submitted more than two days late will not receive feedback and will automatically earn a grade of zero. If your failure to turn your work in on time was due to a University-excused absence, please contact your instructor, and make-up work can be arranged.

Responsible Use of Generative Al

It is <u>University of Maryland Policy</u> and the expectation of MAGE that instructors clearly communicate with their students regarding AI usage in their course. Please use this space to <u>clearly specify</u> your policy on AI use in your class. What (if any) Generative AI use is permitted, and what is not? University policy also dictates that you include mention of university approved tools and provide your students with a link to those tools. Below is a sample AI policy that you can edit for your needs. Consider also, specifying appropriate AI use for each of your assignments in the Academic Integrity Chart located in the Academic Integrity section below.

GENERATIVE AI POLICY: Generative AI tools (e.g., ChatGPT, GitHub Copilot, etc.) are becoming increasingly common in engineering education and in the workplace. In this course, students are expected to use AI technologies ethically and in ways that support learning, uphold academic integrity, and align with course objectives.

Permitted Uses of AI Tools in This Course

Students may use generative AI tools for the following purposes:

- Brainstorming initial ideas or outlining for assignments
- Getting help understanding difficult engineering concepts (e.g., asking for explanations or examples)
- Writing assistance at the sentence level (e.g., grammar or clarity improvements)
- Debugging support in coding assignments, provided students still understand and can explain their code

Prohibited Uses of AI Tools in This Course

Students may not use generative AI tools for:

- Completing graded assignments, problem sets, or projects unless explicitly permitted
- Generating solutions to coding or engineering problems without understanding and verifying the output
- Writing full sections of reports, papers, or lab assignments
- Submitting Al-generated work as their own without proper citation or instructor approval

It is the student's responsibility to make sure any use of AI aligns with the expectations outlined above. Misuse of AI tools may constitute academic dishonesty and will be addressed accordingly (see section on academic integrity, below). Lastly, please become familiar with the <u>University-approved AI tools</u> and university guidelines on <u>responsible</u> AI use. If you are unsure whether a particular use of AI is appropriate, please ask before proceeding.

Academic Integrity

For this course, some of your assignments will be collected via Turnitin on ELMS/Canvas. I have chosen to use this tool because it can help you improve your scholarly writing and help me verify the integrity of student work. For information about Turnitin, how it works, and the feedback reports you may have access to, visit Turnitin Originality Checker for Students

The University's Code of Academic Integrity is designed to ensure that the principles of academic honesty and integrity are upheld. In accordance with this code, the University of Maryland does not tolerate academic dishonesty. Please ensure that you fully understand this code and its implications because all acts of academic dishonesty will be dealt with in accordance with the provisions of this code. All students are expected to adhere to this Code. It is your responsibility to read it and know what it says, so you can start your professional life on the right path. As future professionals, your commitment to high ethical standards and honesty begins with your time at the University of Maryland.

It is important to note that course assistance websites, such as CourseHero, or Al-generated content, are not permitted sources unless the instructor explicitly gives permission. Material taken or copied from these sites can be deemed unauthorized material and a violation of academic integrity. These sites offer information that might be inaccurate or biased, and most importantly, relying on restricted sources will hamper your learning process, particularly the critical thinking steps necessary for college-level assignments.

Additionally, students may naturally choose to use online forums for course-wide discussions (e.g., Group lists or chats) to discuss concepts in the course. However, collaboration on graded assignments is strictly prohibited unless otherwise stated. Examples of prohibited collaboration include: asking classmates for answers on quizzes or exams, asking for access codes to clicker polls, etc. Please visit the Office of Graduate Studies' full list of campus-wide policies and reach out if you have questions.

Finally, on each exam or assignment you must write out and sign the following pledge: "I pledge on my honor that I have not given or received any unauthorized assistance on this exam/assignment."

If you ever feel pressured to comply with someone else's academic integrity violation, please reach out to me straight away. Also, *if you are ever unclear* about acceptable levels of collaboration, *please ask*!

To help you avoid unintentional violations, *the following table* lists levels of collaboration that are acceptable for each graded exercise. Each assignment will contain more specific information regarding acceptable levels of collaboration.

	Open Notes	Use Book	Learn Online	Gather Content with AI	Ask Friends	Work in Groups
Homeworks	•	~	~			

9 Last updated 11/3/2025

In-Class Problem Sets & Quizzes	✓	•	✓		>	~
Final Project	v	•	v	•	>	v
Midterm Exam						

Course Evaluation

Please submit a course evaluation through Student Feedback on Course Experiences in order to help faculty and administrators improve teaching and learning at Maryland. All information submitted to Course Experiences is confidential. Campus will notify you when Student Feedback on Course Experiences is open for you to complete your evaluations at the end of the semester. Please go directly to the <u>Student Feedback on Course Experiences</u> to complete your evaluations. By completing all of your evaluations each semester, you will have the privilege of accessing through Testudo the evaluation reports for the thousands of courses for which 70% or more students submitted their evaluations.

Religious Observance

It is the student's responsibility to inform the instructor of any intended absences for religious observances in advance. Notice should be provided as soon as possible, but no later than the end of the schedule adjustment period.

Copyright Notice

Course materials are copyrighted and may not be reproduced for anything other than personal use without written permission.

Tips for Succeeding in this Course

- 1. **Participate.** I invite you to engage deeply, ask questions, and talk about the course content with your classmates. You can learn a great deal from discussing ideas and perspectives with your peers and professor. Participation can also help you articulate your thoughts and develop critical thinking skills.
- 2. **Manage your time.** Students are often very busy, and I understand that you have obligations outside of this class. However, students do best when they plan adequate time that is devoted to coursework. Block your schedule and set aside plenty of time to complete assignments, including extra time to handle any technology-related problems.
- 3. **Log in regularly.** I recommend that you log in to ELMS-Canvas several times a week to view announcements, discussion posts, and replies to your posts. You may need to log in multiple times a day when group submissions are due.
- 4. **Do not fall behind.** This class moves at a quick pace and each week builds on the previous content. If you feel you are starting to fall behind, check in with the instructor as soon as possible so we can troubleshoot together. It will be hard to keep up with the course content if you fall behind in the pre-work or post-work.
- 5. **Use ELMS-Canvas notification settings.** Pro tip! Canvas ELMS-Canvas can ensure you receive timely notifications in your email or via text. Be sure to enable announcements to be sent instantly or daily.
- 6. **Ask for help if needed.** If you need help with ELMS-Canvas or other technology, IT Support. If you are struggling with a course concept, reach out to me and your classmates for support.

Student Resources and Services

Taking personal responsibility for your learning means acknowledging when your performance does not match your goals and doing something about it. I hope you will come talk to me so that I can help you find the right approach to success in this course, and I encourage you to visit the <u>Counseling Center's Academic Resources</u> to learn more about the wide range of resources available to you. Below are some additional resources and services commonly used by graduate students. For a more comprehensive list, please visit the Graduate School's <u>Campus Resources Page</u>.

Accessibility and Disability Services

The University of Maryland is committed to creating and maintaining a welcoming and inclusive educational, working, and living environment for people of all abilities. The University of Maryland is also committed to the principle that no qualified individual with a disability shall, on the basis of disability, be excluded from participation in or be denied the benefits of the services, programs, or activities of the University, or be subjected to discrimination. The Accessibility & Disability Service (ADS) provides reasonable accommodations to qualified individuals to provide equal access to services, programs and activities. ADS cannot assist retroactively, so it is generally best to request accommodations several weeks before the semester begins or as soon as a disability becomes known. Any student who needs accommodations should contact me as soon as possible so that I have sufficient time to make arrangements.

For assistance in obtaining an accommodation, contact Accessibility and Disability Service at 301-314-7682, or email them at adsfrontdesk@umd.edu. Information about sharing your accommodations with instructors, note taking assistance and more is available from the Counseling Center.

Writing Center

Everyone can use some help sharpening their communication skills (and improving their grade) by visiting <u>The Graduate School's Writing Center</u> and schedule an appointment with them. Additionally, international graduate students may want to take advantage of the Graduate School's free <u>English Editing for International Graduate Students (EEIGS) program</u>.

Health Services

The University offers a variety of physical and mental health services to students. If you are feeling ill or need non-emergency medical attention, please visit the <u>University Health Center</u>.

If you feel it would be helpful to have someone to talk to, visit <u>UMD's Counseling Center</u> or <u>one of the many other</u> mental health resources on campus.

Notice of Mandatory Reporting

Notice of mandatory reporting of sexual assault, sexual harassment, interpersonal violence, and stalking: As a faculty member, I am designated as a "Responsible University Employee," and I must report all disclosures of sexual assault, sexual harassment, interpersonal violence, and stalking to UMD's Title IX Coordinator per University Policy on Sexual Harassment and Other Sexual Misconduct.

If you wish to speak with someone confidentially, please contact one of UMD's confidential resources, such as <u>CARE</u> to <u>Stop Violence</u> (located on the Ground Floor of the Health Center) at 301-741-3442 or the <u>Counseling Center</u> (located at the Shoemaker Building) at 301-314-7651.

You may also seek assistance or supportive measures from UMD's Title IX Coordinator, Angela Nastase, by calling 301-405-1142, or emailing titleIXcoordinator@umd.edu.

To view further information on the above, please visit the <u>Office of Civil Rights and Sexual Misconduct's</u> website at <u>ocrsm.umd.edu</u>.

Basic Needs Security

If you have difficulty affording groceries or accessing sufficient food to eat every day, or lack a safe and stable place to live, please visit <u>UMD's Division of Student Affairs website</u> for information about resources the campus offers you and let me know if I can help in any way.

Veteran Resources

UMD provides some additional supports to our student veterans. You can access those resources at the office of <u>Veteran Student life</u> and the <u>Counseling Center</u>. Veterans and active duty military personnel with special circumstances (e.g., upcoming deployments, drill requirements, disabilities) are welcome and encouraged to communicate these, in advance if possible, to the instructor.