

	2026		1		CSED490H		01
							3-0-3
	:						
	(14:00 15:15) - 2		[109]				
E-Mail	SANGDON@POSTECH.AC.KR		Homepage		HTTPS://SANGDON.GITHUB.IO/		
					054-279-2396		
Office Hours							
As AI advances and is being-practical, its safety and security concerns dramatically emerge. In this class, we learn the art of attacking AI systems along with necessary concepts and tools in AI. In particular, we will learn two core concepts, victim models (e.g., LLMs, VLAs, and Agentic AI) and attack methods (e.g., adversarial examples and jailbreaking) along with core optimization tools (e.g., gradient descent, policy optimization, and prompt tuning with LoRA). At the end of this class, students will have a good understanding of trendy AI models, broad aspects of AI red teaming methods, and necessary AI tools. Note that this course is designed for undergraduates -- graduate students may audit.							
	/						
- Artificial Intelligence							
	가						
					/	/	
- Assignment/Presentation 80% -- three HWs and one final project# - Participation: 20%							
							ISBN
Related references include the following:#							

	2026		1		CSED490H		01
							3-0-3
	:						
	(14:00 15:15) - 2			[109]			
<ul style="list-style-type: none"> - Ian J. Goodfellow et al., "Explaining and Harnessing Adversarial Examples," ICLR '15.# - Ashish Vaswani et al., "Attention Is All You Need," NIPS '17.# - John Schulman et al., "Trust Region Policy Optimization," ICML '15. 							
<p>Week 1: #</p> <ul style="list-style-type: none"> - Introduction to AI Security# <p>Week 2: #</p> <ul style="list-style-type: none"> - Preliminary: Neural Networks / SGD # - Inference-time Attacks: Adversarial Examples / Adversarial Patches / Transfer Attacks# <p>Week 3: #</p> <ul style="list-style-type: none"> - Preliminary: Transformers / LLMs / LCMs / LRM# - Preliminary: RAG# <p>Week 4:#</p> <ul style="list-style-type: none"> - Student Presentation and Discussion on HW 1# <p>Week 5:#</p> <ul style="list-style-type: none"> - Preliminary: Diffusion Models# - Preliminary: Vision-Language-Action Models# <p>Week 6:#</p> <ul style="list-style-type: none"> - Student Presentation and Discussion on HW 2# <p>Week 7:#</p> <ul style="list-style-type: none"> - Preliminary: Optimization for Whitebox Victim Models -- Prompt tuning methods (e.g., LoRA) # - Preliminary: Optimization for Blackbox Victim Models -- Zero-th Order Optimization# <p>Week 8:#</p> <ul style="list-style-type: none"> - Preliminary: Optimization for Blackbox Victim Models -- RL / Policy Optimization# - Inference-time Attacks: Prompt Leaking, Prompt Injection, Jailbreaking# <p>Week 9:#</p> <ul style="list-style-type: none"> - Preliminary: Agentic AI / Tool-calling Agents# - Inference-time Attacks: Current Trends on Red Teaming# 							

	2026		1		CSED490H		01
							3-0-3
	:						
	(14:00 15:15) - 2			[109]			

Week 10:#

- Student Presentation and Discussion on HW 3#

Week 11:#

- Training-set Attacks: membership inference attacks#
- Training-set Attacks: data poisoning attacks#

Week 12:#

- Model Attacks: model extraction attacks#

Week 13:#

- Final Remarks: Overview on defense methods#

Week 14:#

- Student Presentation and Discussion on Final Projects#

Week 15:#

- Student Presentation and Discussion on Final Projects