

2022년도 2학기 전자전기공학특론D(효율적인 인공지능 모델) (EECE695D-01) 강의계획서

1. 수업정보

학수번호	EECE695D	분반	01	학점	3.00
이수구분	전공선택	강좌유형	강의실 강좌	선수과목	
포스테키안 핵심역량	<input type="checkbox"/> 대인관계역량 <input type="checkbox"/> 글로벌시민역량 <input type="checkbox"/> 지식탐구역량 <input type="checkbox"/> 디지털리터러시역량 <input type="checkbox"/> 자기관리역량 <input type="checkbox"/> 창의융합역량				
강의시간	화, 목 / 09:30 ~ 10:45 / 제2공학관 강의실 [102호]			성적취득 구분	G

2. 강의교수 정보

	이름	이재호	학과(전공)	전자전기공학과
	이메일 주소	jaehoklee@postech.ac.kr	Homepage	
	연구실		전화	
	Office Hours	By appointment (Thu 17:00~18:00).		

3. 강의목표

This course gives you an overview on the emerging area of "Efficient ML," i.e., building effective ML models under various resource constraints (e.g., active GPU RAM, total training FLOPs, inference latency). The main focus is on the algorithmic advances for efficient ML, while we still discuss relevant issues on the system/hardware side.

4. 강의선수/수강필수사항

No formal requirements, but I expect you to have some elementary knowledge on machine learning and deep learning.
 If you took EECE454 (or AIGS515/538 or CSED490X/Y), you should be okay.
 Feel free to contact me if you are worried :)

5. 성적평가

Homework Assignments: 30%
 In-class Presentation: 30%
 Final Report: 30%
 Participation: 10%

6. 강의교재

도서명	저자명	출판사	출판년도	ISBN
Lecture notes will be given, but I still recommend you to read the materials below.			0000	

7. 참고문헌 및 자료

Sze, Chen, Yang, and Emer, "Efficient Processing of Deep Neural Networks," Morgan Claypool, 2020.
 Menghani, and Singh, "Efficient Deep Learning," (not-published-yet), 2022.

8. 강의진도계획

W1: Introduction to EfficientML & Basics of Deep Learning
 W2-3: Compute and Memory in Deep Learning, Scaling Law
 W4-5: Network Architectures / NAS
 W6: Weight Initialization, Transfer Learning, and Meta-Learning
 W7: AutoML, Hyperparameter Tuning
 W8: Data and Model Parallelism in Deep Learning
 W9-11: Model Compression
 W12: Efficient Training
 W13: TinyML

W14-15: Student Presentations

W16: Final Report

9. 수업운영

- Final Report: A short survey of a small sub-field of EfficientML, that is most relevant to your research or interest.
- Student Presentation: A 20-minute-long talk about a paper published in NeurIPS, ICML, ICLR or MLSys (2020-2022), or a more recent preprint.
- Homework Assignments: Involves coding and math; get ready for it!
- Participation: Not about attendance (you are all professionals), but about an active discussion.

10. 학습법 소개 및 기타사항

11. 장애학생에 대한 학습지원 사항

- 수강 관련: 문자 통역(청각), 교과목 보조(발달), 노트필기(전 유형) 등
- 시험 관련: 시험시간 연장(필요시 전 유형), 시험지 확대 복사(시각) 등
- 기타 추가 요청사항 발생 시 장애학생지원센터(279-2434)로 요청