

# Jimin Yeom | Curriculum Vitae

✉ sion0107@hanyang.ac.kr • 🌐 jiminyeom.com

• LinkedIn: <https://www.linkedin.com/in/jimin-yeom-a5435b33b/>

## Research Interests

---

My research goal is to build **trustworthy and interpretable AI** by understanding how models learn. Concretely, I study the optimization problem

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \mathcal{L}(f(x + \delta, \theta + \eta), y) \right]$$

which unifies **adversarial training** (perturbations  $\delta$  in input space) and the study of **training dynamics** (perturbations  $\eta$  in parameter space). I view this joint objective as a lens through which robustness, generalization, and interpretability can be analyzed together, rather than as separate phenomena.

**Trustworthy AI:** *Adversarial Robustness* – Robustness–Accuracy Trade-off, Catastrophic Overfitting; *LLM Jailbreaking* – Certifiable Defense.

**Interpretable AI:** *In-Context Learning*, *Linear Transformer*, *Model Unlearning*, *Reasoning*.

## Education

---

**Hanyang University**

*M.S. in Computer Science*

Department of Computer Science & Software

Advisor: *Sungyoon Lee*

**Seoul, Republic of Korea**

*Sep. 2025 – Aug. 2027 (expected)*

**Hanyang University**

*B.S. in Computer Science*

Department of Computer Science & Software

**Seoul, Republic of Korea**

*Mar. 2018 – Aug. 2025*

## Publications

---

\* denotes equal contribution. All papers listed below are currently **under review**.

**[Paper 1]**

**Jimin Yeom**, Jonghyun Hong, *Sungyoon Lee*.

*Keywords:* AI Safety, LLM Jailbreaking, Optimization.

*Under Review*

**[Paper 2]**

**Jimin Yeom**, *Sungyoon Lee*.

*Keywords:* AI Safety, Adversarial Robustness, Adversarial Training.

*Under Review*

**[Paper 3]**

Subin Jang, **Jimin Yeom**, *Sungyoon Lee*.

*Keywords:* AI Safety, LLM Jailbreaking, Certified Defense.

*Under Review*

## Ongoing Work

---

*Edge of Stability in Adversarial Training*

Investigating how the Edge-of-Stability regime manifests under adversarial training and how it interacts with catastrophic overfitting.

*2025 – present*

*Reward Hacking*

Studying the mechanisms behind reward hacking in alignment-tuned models, with the goal of mitigating it via training-dynamics-aware objectives.

*2025 – present*