# Physical AI security primer

Shin, Jongho/LG Electronics

# Who am I

- Research Fellow at LG Electronics

- Affiliate professor in Kookmin university

- Cyber security expert

  - W/ long history of offensive security

- Research focuses

  - AI security

  - Privacy Enhancing Technology

# Physical AI

## Meta's Yann LeCun to Launch Physical A.I. Startup After Declaring LLMs a 'Dead End'

The pioneering A.I. researcher is betting on a new paradigm that teaches machines to understand the physical world, not just language.

By Alexandra Tremayne-Pengelly • 11/11/25 2:09pm

## Nvidia CEO Jensen Huang Predicts the Next Big Thing After 'Agentic A.I.'

"Now is the beginning of the agentic A.I. era...then there's physical A.I. after that."

By Alexandra Tremayne-Pengelly • 02/27/25 3:04pm

"The next wave is already happening... Robotics, which has been enabled by physical AI, AI that understands the physical world." - Jensen Huang

# Home robots are coming



SCALING HELIX
LAUNDRY

# Moving surveillance camera



We hacked a robot vacuum —
and could watch live through its
camera

https://www.abc.net.au/news/2024-10-04/robot-vacuum-hacked-photos-camera-audio/104414020

# But

- Physical AI security is more than just device hackings.

# Evolution of AI

## Gen AI

Generates new content

Uses deep learning models

Trained on large datasets

## AI agent

Perform tasks for users

May incorporate both types

Interactive & collaborative

## Physical AI

Interacts with the physical world

Uses sensors & actuators

Converts decisions into physical
actions

# Part 1 — AI Security

- Input → Model → Output → (post processing)



**GenAI Framework**

Prompt

Prompt Engineering

Inputs

GenAI

Fine-Tuning

Reinforcement Learning

Outputs

Governance

Source: Gartner
© 2024 Gartner, Inc. and/or its affiliates. All rights reserved. 2948768

**Gartner.**

# Part 1 — AI Security

- Attack vectors

  - Prompt Injection

  - Jailbreak

  - Hallucination

  - Model Extraction

  - RAG Manipulation

  - MLOps Pipeline Attack

- These are **bounded** and **predictable**.

# Part 2 — Agent Security

- AI Agent: reason, plan, and act

- Agents can autonomously:
  - 1. Use tools
  - 2. Call external APIs
  - 3. Read/write files
  - 4. Navigate the web
  - 5. Execute OS-level actions
  - 6. Collaborate with other agents

FIGURE 1: | The core components of an AI agent



Source: World Economic Forum

# Part 2 — Agent Security

- New attack vectors
  - Tool Injection
  - Action Hijacking
  - Delegated Misbehavior
  - Multi-agent Escalation
  - State Manipulation
  - Goal Drift / Value Hijack

# Prompt injection example

https://www.blazeinfosec.com/post/llm-pentest-agent-hacking/

# Part 3 — Physical AI Security

- Physical AI: perceive the physical world, reason, plan, and act within it

- Example

  - Robot vacuums

  - Humanoids

  - Industrial robots

  - Self-driving car/drone

  - Smart home

# Part 3 — Physical AI Security

- VLA(Vision Language Action) model

  - Sensor spoofing

  - Adversarial audio input

  - Adversarial vision patch

  - Model poisoning



https://openvla.github.io/

# When physical security fails

# Complexity and impact explosion

- Innate complexity, attack surface, interactions

- AI Security: Protects **models**

- Agent Security: Protects **actions**

- Physical AI Security: Protects the **real world**

# Stage 1 — AI Models: Linear Complexity

- Security issues at the model level grow linearly:

  - Larger models → more parameters

  - Longer context → more prompt surfaces

  - Added RAG pipelines → more data pathways

# Stage 2 — Agentic AI: The Inflection Point of Non-Determinism

- Agents can autonomously do various things.

- This introduces **state changes**, meaning the system is no longer static or predictable.

  - Attack targets the Action Layer, not the model.

  - Vulnerability comes from inter-agent dynamics, not a single model.

# Stage 3 — Combinatorial Explosion

- With agents, actions no longer follow linear sequences.

- They form graph-shaped decision trees across:

  - Tools

  - States

  - External APIs

  - Histories

  - Environmental input

- This creates a near-infinite combinatorial space

# Stage 4 — Introduction of Environment: Infinite Variables

- Agents connected to the physical world must interpret:

  - Light

  - Sound

  - Temperature

  - Obstacles

  - Human motion

  - Random noise

  - Sensor uncertainty

- The environment introduces infinite, uncontrollable variables.

# Stage 5 — Physical AI: Security Meets Safety
## Security becomes a multi-disciplinary problem.

- Once AI actions influence physical actuators, security failures become safety hazards.

- Examples:

  - Smart oven sets incorrect temperature

  - Home robot moves aggressively toward an object

  - Industrial robot arm miscalculates trajectory

  - Vehicle AI misinterprets traffic conditions

  - HVAC/air-quality AI reacts to spoofed sensors

# The security challenge grows exponentially

- Expansion of the protection scope

  - Model → Behavior → Environment → Physical impact

    - transforming digital errors into real-world harm

- Attack surface transition

  - Linear → Graph-shaped → Infinite

- State changes and non-determinism

  - making static verification impossible

- Environmental variables

  - introducing uncontrolled real-world noise

# So what should we do?

- We need **dynamic** risk assessment; environment itself can be attack surface

- Security must handle **unbounded state spaces**.

- AI red teaming

  - We need more of likeminded people who can see the things in different perspectives.

    - More complexity requires more insights

# Physical AI security initiative



"피지컬AI 오작동 방지"...에임인텔리전스, LG전자와
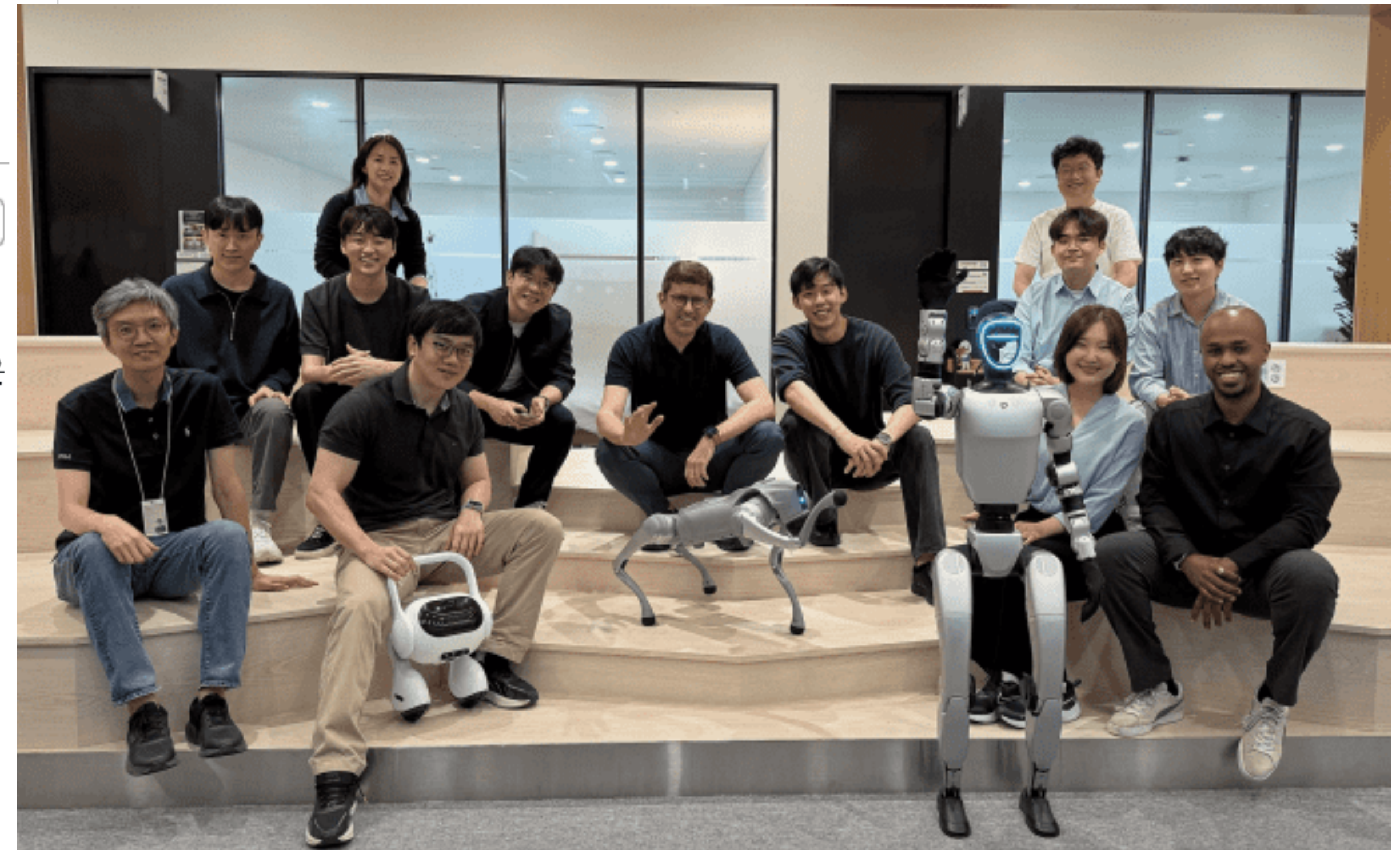공동 개발 추진

| 미국 로봇OS 기업 오픈마인드도 참여..."외부 공격까지 차단하는 체계 구축"

컴퓨팅 | 입력 :2025/10/14 22:17

방은주 기자 | ✉  Ⓝ 기자 페이지 구독  📄 기자의 다른기사 보기   f 𝕏 in 가+ 가- 🖨

AI 보안기업 에임인텔리전스(AIM Intelligence, 대표 유상윤)는 미국 로봇OS 기업 오픈마인드(OpenMind),
LG전자와 함께 '피지컬 AI 안전 레이어(Physical AI Safety Layer)'를 공동 개발한다고 14일 밝혔다. 이번 협력은
물리적 환경에서 작동하는 AI, 즉 '피지컬 AI(Physical AI)'의 오작동과 위험을 방지하기 위한 것이다.

https://zdnet.co.kr/view/?no=20251014221730

# Security community

- AI is shifting from generating text → taking actions → affecting the physical world.

  - Security must evolve with it.

- Prompt Zero

  - https://discord.gg/5TaxVHVP86

# Thank you

- Email
  - jongho0x80@gmail.com
- Linkedin
  - https://kr.linkedin.com/in/shinjongho