

Lecture 1: Introduction to Generative and Explainable AI

Generative AI

<https://www.youtube.com/watch?v=NRmAXDWJVnU>

Agentic AI

https://www.youtube.com/watch?v=15_pppse4fY

Generative AI is a branch of artificial intelligence focused on **creating new content**—such as text, images, audio, video, or even code—based on patterns it has learned from existing data.

Unlike traditional AI systems, which are mainly designed for classification, prediction, or decision-making, generative AI models can **produce original outputs** that resemble human-created work.

Key Points:

- **Definition:** Generative AI refers to machine learning models (often deep learning) that can generate new data instances similar to the training data.
- **How it works:**
 - It uses large datasets to learn underlying patterns, structures, and relationships.
 - Techniques like *Generative Adversarial Networks (GANs)*, *Variational Autoencoders (VAEs)*, and *Transformer-based models* (e.g., GPT, Stable Diffusion, MidJourney) are common.
- **Examples:**
 - ChatGPT → generates human-like text.
 - DALL·E, Stable Diffusion → generate images from text prompts.
 - MusicLM → creates music.
 - GitHub Copilot → generates computer code.
- **Applications:**
 - Content creation (articles, art, design).
 - Healthcare (drug discovery, medical imaging synthesis).

- Education (personalized learning content).
- Business (marketing copy, product design).

Explainable AI (XAI) refers to methods and techniques that make the results of AI systems **transparent, understandable, and interpretable** to humans.

Since many modern AI models (especially deep learning ones) act like a **"black box"**—they make predictions without showing how or why—XAI is about **opening that box** and giving insights into the decision-making process.

◆ Why is XAI important?

- **Trust** → Users, doctors, regulators, or businesses need to trust AI's output.
 - **Accountability** → If something goes wrong (e.g., medical misdiagnosis, credit rejection), we must know *why*.
 - **Bias detection** → Helps spot unfair or discriminatory patterns.
 - **Regulatory compliance** → Laws like the EU AI Act and GDPR demand transparency in automated decision-making.
-

◆ How it works


XAI provides explanations in forms humans can understand, such as:

1. **Feature importance** – showing which input features most influenced the decision.
 2. **Visualization** – heatmaps in images (e.g., showing which part of an X-ray led to a diagnosis).
 3. **Rule-based explanations** – converting model behavior into “if-then” rules.
 4. **Counterfactual explanations** – “If X had been different, the outcome would have changed.”
-

◆ Examples

- **Healthcare:** An AI says “tumor detected.” XAI highlights the tumor region in the scan to justify.
- **Finance:** Loan denied → XAI shows the key factors (e.g., low income, poor credit history).
- **Self-driving cars:** Explains why the car braked suddenly (e.g., pedestrian detected).

◆ Difference between Generative AI and Explainable AI

Aspect	Generative AI	Explainable AI	
Goal	Create new content (text, image, audio, etc.)	Make AI decisions understandable to humans	
Nature	Creative, produces novel outputs	Analytical, provides reasoning behind outputs	
Examples	ChatGPT, DALL·E, Stable Diffusion	SHAP, LIME, Grad-CAM	
Challenge	Risk of hallucination, bias in generated data	Complexity of explaining deep models	
Applications	Art, writing, design, drug discovery	Healthcare, finance, law, autonomous systems	