

AI Demystified for CXO

Jaideep Ganguly, Sc.D. (MIT)

<https://jganguly.github.io>

23 Jan 2026

About the Speaker

- 
- 35 years of industry experience in India and USA.
 - Technical Fellow & Consultant
 - ▶ Quantum Computing Algorithms in Reinforced Learning
 - ▶ AI in Medicine
 - Amazon, Technology Leader
 - Microsoft, Technology Leader
 - Massachusetts Institute of Technology (MIT),
Doctor of Science, Master of Science
 - Indian Institute of Technology, Kharagpur,
B.Tech (Hons.)

Outline

1. What is AI?
2. History & Evolution of AI
3. AI, ML, Deep Learning, LLM, Gen AI
4. Why AI Matters Now
5. The AI Stack: A CXO View
6. Where AI Creates Business Value
7. AI Use Cases: What Works, What Doesn't
8. Core AI Technologies
9. Deep Neural Network
10. Activation
11. Back Propagation
12. LLM - Large Language Model
13. Attention is all we need
14. Generative Pre-trained Transformer
15. Input Representation
16. Input Representation: Math Formulation
17. Encoder
18. Encoder: Math Formulation
19. Decoder
20. Decoder: Math Formulation
21. Why is Gen AI / LLM a revolution?
22. What can the LLMs do?
23. Some Large Language Models
24. Closing Thoughts & Questions

What is Artificial Intelligence (AI)?

- **Artificial Intelligence (AI)** refers to systems that perform tasks which typically require human intelligence. These tasks include:
 - ▶ Perception (vision, speech)
 - ▶ Reasoning and decision-making
 - ▶ Learning from data
 - ▶ Problem solving and planning
- Core paradigms of AI
 - ▶ **Symbolic AI:** Rules, logic, knowledge graphs
 - ▶ **Statistical AI:** Probabilities, inference, optimization
 - ▶ **Neural AI:** Learning representations from data
- Modern AI systems are primarily based on:
 - ▶ statistical learning
 - ▶ optimization
 - ▶ large-scale computation
- AI does *not* imply consciousness or general intelligence; most systems are task-specific.

History & Evolution of AI

● Foundations (1950s–1960s)

- ▶ Alan Turing (University of Manchester): Turing Test, computability
- ▶ John McCarthy (MIT / Stanford): coined the term “Artificial Intelligence”
- ▶ Marvin Minsky (MIT): symbolic AI, early cognitive architectures

● Expert Systems Era (1970s–1980s)

- ▶ Edward Feigenbaum (Stanford University): knowledge-based systems
- ▶ MYCIN and DENDRAL expert systems
- ▶ Patrick Winston (MIT): knowledge representation, AI pedagogy

● Statistical Machine Learning (1990s–2000s)

- ▶ Vladimir Vapnik (AT&T Bell Labs): Support Vector Machines
- ▶ Judea Pearl (UCLA): probabilistic reasoning, Bayesian networks
- ▶ Growing influence of MIT in robotics and learning-based AI

● Deep Learning Revolution (2010s)

- ▶ Geoffrey Hinton (University of Toronto / Google): Deep Neural Networks
- ▶ Yann LeCun (AT&T Bell Labs / Meta): convolutional neural networks
- ▶ Yoshua Bengio (Université de Montréal): representation learning
- ▶ MIT CSAIL: advances in robotics, perception, and human-centered AI

● Foundation Models (2020s–Present)

- ▶ Transformer architecture (Google Brain)
- ▶ Large-scale models developed at OpenAI, DeepMind, Meta
- ▶ Ongoing contributions from MIT in AI safety, systems, and applied AI

AI, ML, Deep Learning, LLM, Generative AI

- **Artificial Intelligence (AI)**

- ▶ Broad goal: systems that exhibit human-like intelligence
- ▶ Includes reasoning, planning, perception, and decision-making

- **Machine Learning (ML)**

- ▶ Subset of AI focused on learning patterns from data. Improves performance without explicit rule-based programming
- ▶ Types of ML
 - ★ **Supervised Learning:** Labeled data (prediction, classification)
 - ★ **Unsupervised Learning:** Structure discovery (clustering)
 - ★ **Reinforcement Learning:** Learning via feedback and rewards

- **Deep Learning (DL)**

- ▶ Subset of ML using multi-layer neural networks
- ▶ Excels at vision, speech, language, and complex representations

- **LLMs (Large Language Models)** are a specific type of model that generate text and language-based outputs. All LLMs are Generative AI

- **Generative AI (GenAI)**

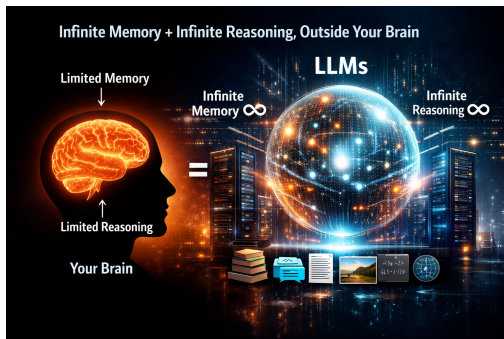
- ▶ Subset of DL focused on generating new content; Produces text, images, code, audio, and simulations.
- ▶ Not all Generative AI are LLMs.

Each layer is a specialization of the one above it

Why AI Matters Now

- AI is not new — but **three forces converged** in the last decade:
 - ▶ Explosion of digital data
 - ▶ Massive, affordable compute (cloud, GPUs)
 - ▶ Breakthrough algorithms (deep learning, transformers)
- AI has crossed a threshold:
 - ▶ From research labs to real business deployment
 - ▶ From niche tools to enterprise platforms
- Generative AI made AI *visible* and accessible to everyone

Key message: AI is now a competitive necessity, not an experiment.



The AI Stack: A CXO View

- **Data Layer**

- ▶ Enterprise data, quality, ownership, governance

- **Model Layer**

- ▶ Classical ML, deep learning, foundation models

- **Application Layer**

- ▶ Copilots, automation, decision-support systems

- **Infrastructure Layer**

- ▶ Cloud, on-prem, security, cost management

Key message: AI value depends on the *entire stack*, not just models.

Where AI Creates Business Value

● Revenue Growth

- ▶ Personalization
- ▶ pricing optimization
- ▶ sales effectiveness

● Cost Reduction

- ▶ Automation
- ▶ Process efficiency
- ▶ Reduced manual effort

● Risk Reduction

- ▶ Fraud detection
- ▶ Compliance
- ▶ Anomaly detection

● Decision Quality & Speed

- ▶ Forecasting
- ▶ Scenario analysis
- ▶ Recommendations

Key message: AI is a business lever, not an IT project.

AI Use Cases: What Works, What Doesn't

● High-ROI Use Cases

- ▶ Demand forecasting
- ▶ Recommendations and personalization
- ▶ Document and knowledge automation
- ▶ Customer support copilots

● High-Risk / Overhyped Use Cases

- ▶ Fully autonomous decision-making
- ▶ Replacing judgment-heavy roles
- ▶ AI without clean or owned data

● Why Pilots Fail

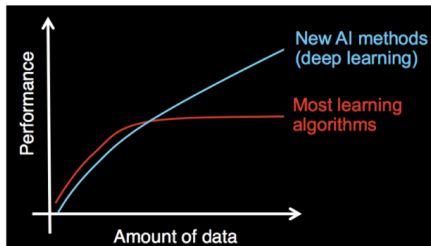
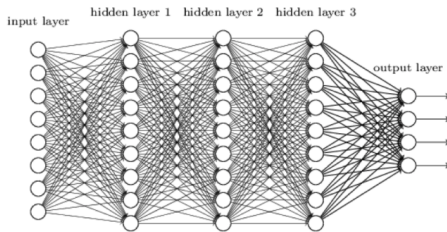
- ▶ Unclear business ownership
- ▶ Poor data foundations
- ▶ No success metrics

Key message: Successful AI starts with the right problem, not the right model.

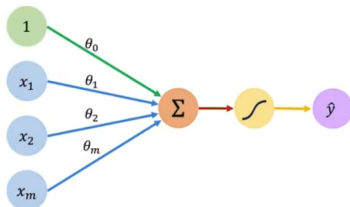
Core AI Technologies

- 1 Deep Neural Network
- 2 Generative AI (Text, Image, Audio, Video)
- 3 Large Language Models (LLMs)
- 4 Multimodal AI
- 5 Reinforcement Learning (RL)
- 6 AI Agents & Autonomy
- 7 Explainable AI (XAI)

Deep Neural Network



Activation



Inputs Weights Sum Non-Linearity Output

Output

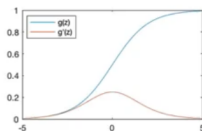
Linear combination of inputs

$$\hat{y} = g \left(\theta_0 + \sum_{i=1}^m x_i \theta_i \right)$$

Non-linear activation function

Bias

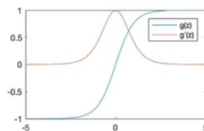
Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

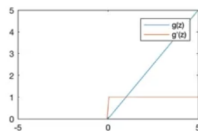
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

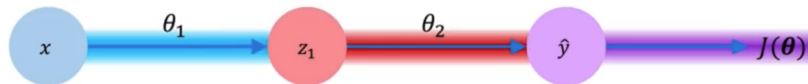
Rectified Linear Unit (ReLU)



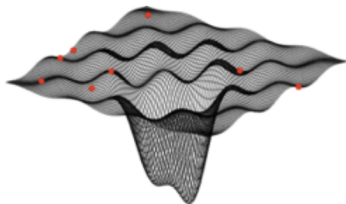
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

Back Propagation



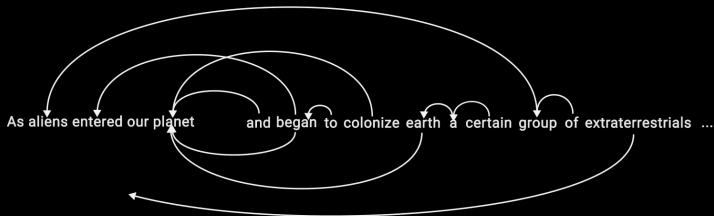
$$\frac{\partial J(\theta)}{\partial \theta_1} = \underbrace{\frac{\partial J(\theta)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial \theta_1}}_{\text{blue}}$$



LLM - Large Language Model

- 1 **LLM:** An LLM learns patterns from massive text datasets and uses them to generate or interpret language. It relies on a **transformer architecture with self-attention**.
- 2 **Transformer:** A type of **neural network** used by modern AI models (like GPT, BERT, Claude). Its core innovation is the self-attention mechanism.
- 3 **Self-attention:** The model understands the **relationships between words** in a sentence all **at the same time**, rather than sequentially.
- 4 **Training:** **Processes billions of examples and adjusts billions of parameters** to predict the next word with increasing accuracy, gradually learning grammar, facts, and reasoning patterns.
- 5 **Doesn't think like a human:** When prompted, it does not reason or understand; it simply predicts the most likely next token based on learned patterns. Fine-tuning and reinforcement learning further refine its behavior. Ultimately, an LLM is a **statistical engine** that produces human-like text by modeling deep linguistic patterns, **not by truly understanding the world**.

Attention is all we need



Attention is the core idea that lets a Transformer decide which parts of the input matter most to each other, regardless of distance.

Instead of reading a sequence step-by-step, the Transformer looks at all tokens at once and computes weighted relationships between them.

LLMs do their job using GPT



Input Representation

Embedding and Positional Encoding

- 1 The model begins with a sequence of discrete input tokens representing the source data.
- 2 Each token is transformed into a fixed-length continuous vector using an embedding layer. These vectors capture semantic information in a high-dimensional space.
- 3 All token embeddings share the same dimensionality, known as the model dimension.
- 4 Because embeddings alone do not contain information about word order, an explicit positional encoding is generated for each token position.
- 5 The positional encoding assigns a unique pattern to each position, allowing the model to distinguish between tokens appearing at different locations in the sequence.
- 6 The final input to the encoder is obtained by adding the positional encoding to the corresponding token embeddings.
- 7 This combined representation contains both semantic meaning and positional information and is passed to the encoder stack.

Input Representation: Math Formulation

Embedding + Positional Encoding

$$X_{\text{emb}} = \text{Embed}(x_{1:N_{\text{src}}})$$
$$X_{\text{emb}} = \begin{bmatrix} \text{token } x_1 \\ \text{token } x_2 \\ \vdots \\ \text{token } x_{N_{\text{src}}} \end{bmatrix} \longrightarrow \begin{bmatrix} \mathbf{e}_1 \in \mathbb{R}^{d_{\text{model}}} \\ \mathbf{e}_2 \in \mathbb{R}^{d_{\text{model}}} \\ \vdots \\ \mathbf{e}_{N_{\text{src}}} \in \mathbb{R}^{d_{\text{model}}} \end{bmatrix}$$

Each row of X_{emb} represents the embedding vector of a single input token.

The embedding dimension is the length of the vector used to represent each token (or item) in a continuous space and is typically $d_{\text{model}} \in \{512, 768, 1024, \dots\}$.

$$P_{\text{pos}, 2i} = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \quad i \text{ even}$$

$$P_{\text{pos}, 2i+1} = \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \quad i \text{ odd}$$

$$P_{\text{src}} \in \mathbb{R}^{N_{\text{src}} \times d_{\text{model}}}$$

$$X = X_{\text{emb}} + P_{\text{src}} \in \mathbb{R}^{N_{\text{src}} \times d_{\text{model}}}$$

Encoder

- 1 The encoder processes the entire input sequence simultaneously using self-attention.
- 2 For each token, the model computes multiple attention heads that determine how strongly the token should attend to every other token in the sequence.
- 3 Each attention head focuses on a different type of relationship, such as syntax, semantics, or long-range dependencies.
- 4 The outputs of all attention heads are combined to form a single representation for each token.
- 5 A residual connection and layer normalization are applied to stabilize training and preserve information from earlier layers.
- 6 A position-wise feedforward network is then applied independently to each token representation, introducing nonlinearity and additional expressive power.
- 7 A second residual connection and normalization produce the final encoder output.
- 8 The resulting representations encode the full source sequence and are passed to the decoder as contextual memory.

Encoder: Math Formulation

$$X = X_{\text{emb}} + P_{\text{src}} \in \mathbb{R}^{N_{\text{src}} \times d_{\text{model}}}$$

$$Q^{(i)} = XW_{Q,\text{enc}}^{(i)}, \quad K^{(i)} = XW_{K,\text{enc}}^{(i)}, \quad V^{(i)} = XW_{V,\text{enc}}^{(i)}$$

$$\text{Attn_Head}_i = \text{softmax}\left(\frac{Q^{(i)}(K^{(i)})^\top}{\sqrt{d_k}}\right) V^{(i)}, \quad i = 1, \dots, h$$

$$O = \text{Concat}(\text{Attn_Head}_1, \dots, \text{Attn_Head}_h) W_{O,\text{enc}}$$

$$\mu = \frac{1}{d_{\text{model}}} \sum_{j=1}^{d_{\text{model}}} x_j, \quad \sigma = \sqrt{\frac{1}{d_{\text{model}}} \sum_{j=1}^{d_{\text{model}}} (x_j - \mu)^2 + \epsilon}, \quad \text{LN}(x) = \gamma \odot \frac{x - \mu}{\sigma} + \beta$$

$$x \in \mathbb{R}^{d_{\text{model}}} \text{ (one token vector)}$$

$$X' = \text{LN}(X + O),$$

$$F = \max(0, X'W_{1,\text{enc}} + b_{1,\text{enc}}) W_{2,\text{enc}} + b_{2,\text{enc}},$$

$$E = \text{LN}(X' + F)$$

E = Output of Encoder Block

Decoder

- 1 The decoder operates on the target sequence, which is generated one token at a time.
- 2 Target tokens are first embedded and combined with positional encodings, similar to the encoder input.
- 3 A masked self-attention mechanism ensures that each position can only attend to previously generated tokens, preventing access to future information.
- 4 The decoder then performs cross-attention, allowing each target position to attend to the encoded representations of the source sequence.
- 5 This cross-attention step aligns the target tokens with relevant parts of the input.
- 6 A feedforward network further refines each token representation.
- 7 Residual connections and layer normalization are applied after each major sublayer.
- 8 The final decoder output is projected into the vocabulary space to produce a probability distribution over the next possible token.
- 9 The token with the highest probability is selected or sampled to continue generation.

Decoder: Math Formulation

$$D_{\text{emb}} = \text{Embed}(y_{1:N_{\text{tgt}}}), \quad D = D_{\text{emb}} + P_{\text{tgt}} \in \mathbb{R}^{N_{\text{tgt}} \times d_{\text{model}}}$$

$$Q^{(i)} = DW_{Q,\text{dec}}^{(i)}, \quad K^{(i)} = DW_{K,\text{dec}}^{(i)}, \quad V^{(i)} = DW_{V,\text{dec}}^{(i)}$$

$$\text{Mask_Attn_Head}_i = \text{softmax} \left(\frac{Q^{(i)}(K^{(i)})^\top}{\sqrt{d_k}} + M \right) V^{(i)}, \quad i = 1, \dots, h \quad M_{ij} = \begin{cases} 0 & j \leq i \\ -\infty & j > i \end{cases}$$

$$O_1 = \text{Concat}(\text{Mask_Attn_Head}_1, \dots, \text{Mask_Attn_Head}_h) W_{O,\text{dec}}$$

$$D' = \text{LN}(D + O_1)$$

$$Q^{(i)} = D' W_{Q,\text{cross}}^{(i)}, \quad K^{(i)} = E W_{K,\text{cross}}^{(i)}, \quad V^{(i)} = E W_{V,\text{cross}}^{(i)}$$

$$\text{Cross_Attn_Head}_i = \text{softmax} \left(\frac{Q^{(i)}(K^{(i)})^\top}{\sqrt{d_k}} \right) V^{(i)}$$

$$O_2 = \text{Concat}(\text{Cross_Attn_Head}_1, \dots, \text{Cross_Attn_Head}_h) W_{O,\text{cross}}$$

$$D'' = \text{LN}(D' + O_2),$$

$$F = \max(0, D'' W_{1,\text{dec}} + b_{1,\text{dec}}) W_{2,\text{dec}} + b_{2,\text{dec}},$$

$$Y = \text{LN}(D'' + F)$$

Y = Output of Decoder Block,

$$Z = Y W_{\text{vocab}} + b_{\text{vocab}},$$

$$p(y_t) = \text{softmax}(Z_t)$$

Why is Gen AI / LLM a revolution?

- 1 **New capability:** Machines can now **generate** text, code, images, insights, and decisions—beyond traditional automation.
- 2 **Natural language interface:** Business users can interact with systems using **simple** language, removing technical barriers.
- 3 **Reasoning at scale:** LLMs can analyse documents, summarise, compare, extract insights, and **support decision-making**.
- 4 **Productivity multiplier:** **Dramatically accelerates** work in research, coding, analytics, compliance, design, and communication.
- 5 **Foundation for autonomous workflows:** Enables AI agents that can plan, decide, and act with **minimal human input**.
- 6 **Democratisation of capability:** Advanced AI becomes accessible to **every employee**, not just specialists.

What can the LLMs do?

① Capabilities

- ▶ **Language Understanding & Generation:** Conversation, Q&A, summarisation, translation, rewriting, content creation.
- ▶ **Reasoning & Problem-Solving:** Logical reasoning, planning, multi-step tasks, explanations.
- ▶ **Knowledge & Information:** Concept explanation, knowledge recall, information extraction.
- ▶ **Coding & Technical Tasks:** Code generation, debugging, optimisation, documentation.
- ▶ **Data & Document Intelligence:** Table interpretation, document classification, entity extraction.
- ▶ **Multimodal Capabilities:** Image understanding, speech processing, image/audio generation (model-dependent).
- ▶ **RAG (Retrieval-Augmented Generation):** Search, grounding answers, reducing hallucinations.
- ▶ **Agentic Behaviour:** Tool use, workflow automation, copilots.
- ▶ **Personalisation & Recommendations:** User-specific responses, adaptive learning.

② Limitations: Not infallible

- ▶ May produce incorrect or fabricated information (“hallucinations”)
- ▶ Outputs are based on learned patterns, not real understanding
- ▶ No built-in fact checking or awareness of truth

③ Bottom Line: Very powerful, but must be used with human verification.

Some Large Language Models

Overview

- 1 Foundation models: Gemini, GPT, Claude, Llama
- 2 Multimodal AI: text, image, audio, video
- 3 Agentic systems: planning, reasoning, tool use

Model Lineage & Architecture

- 1 **GPT (OpenAI)** – Dense/hybrid Transformer; natively multimodal (text–image–audio); strong reasoning & tool integration.
- 2 **Gemini (Google)** – Fully multimodal from inception; strong video/audio understanding; hierarchical attention design.
- 3 **Claude (Anthropic)** – Sparse/MoE-enhanced Transformer; long-context specialist; trained with Constitutional AI.
- 4 **Llama (Meta)** – Efficient decoder-only Transformer; open-source; optimized for local, on-device & edge inference.

Key Technical Distinctions

- 1 **Multimodality:** GPT & Gemini are native; Claude uses encoders + fusion; Llama relies on adapters.
- 2 **Training Philosophy:** GPT uses synthetic data + tool traces; Claude uses Constitutional AI; Gemini leverages large multimodal corpora; Llama emphasizes openness & efficiency.
- 3 **Context Windows:** Claude & Gemini lead (1M+ tokens); GPT supports 200k; Llama smaller but extendable.

Closing Thoughts & Questions

*AI is not just tools or code,
but how we choose each path as we go.
It shifts how work and worth align,
and redraws value over time.*

Thank you!