

COMPREHENSIVE REPORT ON THE FUNDAMENTALS OF GENERATIVE AI AND LLMs

1. Foundational Concepts of Gen AI

➤ What Gen AI really is

Generative AI is a class of machine learning systems that learns the statistical structure of data and generate new data samples that resemble it.

It does not:

- Think
- Understand meaning
- Reason like human
- Know facts

It does:

- Learn probability model
- Predict likely continuations
- Recombine patterns seen during training

At its core, Gen AI answers;

"Given what I've seen before, what is the most probable next output?"

➤ Probability modeling:

Generative AI is built on probabilistic modeling

- TEXT → token probability sequences
- IMAGE → pixel or latent distributions
- AUDIO → waveform / spectral distributions

Everything is statistics. No exception.

Generative vs Discriminative models:

The distinction is non-negotiable

MODEL TYPE	What it learns	Example
Discriminative	$P(y)$	(x)
Generative	$P(x)$ or $P(x, y)$	Text / Image generation

Generative models learn how data is formed, not just how to classify it.

Key Limitations:

These are not bugs. That would be

- No factual grounding by default
- No causal understanding
- Hallucinations are inevitable
- Confidence \neq correctness

Anyone claiming "true understanding" is

lying or ignorant.

example: $x \rightarrow y$

bad facts \rightarrow bad output

2. Major AI tools in 2024 (Generative AI Ecosystem)

➤ TEXT & LANGUAGE:

- ChatGPT
- Gemini
- Claude
- LLaMA-based tools
- Perplexity AI
- OpenAI
- Google
- Anthropic
- Meta ecosystem
- Retrieval augmented

➤ IMAGE GENERATION:

- DALL·E
- Midjourney
- Stable Diffusion
- Adobe Firefly

➤ VIDEO GENERATION:

- Sora
- Runway
- Pika Labs

➤ AUDIO & SPEECH:

- ElevenLabs
- Whisper
- Play.ht

➤ CODE + DEVELOPMENT:

- GitHub Copilot
- Cursor
- CodeWhisperer
- Replit AI

➤ PRODUCTIVITY + AUTOMATION:

- Notion AI
- Microsoft Copilot
- Zapier AI
- AutoGPT-style agents

3. Transformer Architecture in Generative AI and Its applications

"Why Transformers Matter?"

Before transformers:

- Sequential processing
- Poor long-range memory
- Slow training

Transformers fixed all of that.

➤ CORE COMPONENTS:

<1> Token embeddings:

converts discrete tokens into dense words/vectors.

<1> Self-Attention Mechanism:

This is the breakthrough.

It allows the model to:

- Weigh relevance of all tokens simultaneously,
- Capture long-range dependencies.
- Process sequences in parallel

FORMALLY:

$$\text{Attention } (\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d})\mathbf{V}$$

No magic. Just matrix math.

<III> Multi-Head Attention:

Multiple attention heads learn different relationships simultaneously.

<IV> Feed-Forward Networks:

Non-linear transformations applied to each token.

<V> Positional Encoding:

Since transformers don't understand order inherently, position information is injected manually.

► Why transformers dominate Generative AI:

- * Scalable
- * Parallelizable
- * Extremely expressive
- * Data-efficient at scale

That's why every serious LLM uses transformers.

► APPLICATIONS:

- * Language modelling
- * Machine translation
- * Code generation
- * Image generation
- * Protein folding
- * Multimodel systems

4. Impact of scaling in Gen AI and LLMs

➤ Scaling Laws:

Empirically observed:

- Larger models + more data + more compute
= better performance.
- Improvements follow predictable power laws

Performance doesn't improve because models get "smarter".

It improves because they get better at statistical approximation.

➤ Emergent abilities:

At certain scales, models suddenly exhibit:

- Few-shot learning
- Tool usage
- Reasoning-like behaviour

These are emergent patterns, not planned capabilities.

> Costs of scaling:

Scaling isn't free

- Massive compute requirements
- Energy consumption
- Centralization of power
- Environmental impact

This is why only a few organizations can train frontier models.

> Diminishing returns:

Each scale increases costs exponentially more for incremental gains.

We're already seeing:

- Plateauing improvements
- Focus shifting to efficiency and alignment.

5. What is a Large Language Model (LLM) and how is it built?

A Large Language Model (LLM) is a deep learning model based on the Transformer architecture that is trained on massive text datasets to predict the next token in a sequence.

By learning statistical patterns in language, LLMs can generate coherent text, answer questions, write code, and perform other language related tasks. Their apparent reasoning ability is an emergent result of scale, not true understanding.

➤ How an LLM is built

1. Data Collection: Large-scale text data is gathered from sources such as books, websites, articles, and code repositories, then cleaned and filtered.

2. Tokenization: Text is broken into smaller units called tokens (subwords or characters) that the model can process numerically.

3. Pretraining: The model is trained using a self-supervised objective - predicting the next token based on previous tokens - across billions or trillions of examples.

4. Fine tuning: The pretrained model is refined using supervised learning and techniques like Reinforcement Learning from Human Feedback (RLHF) to improve instruction-following and safety.

5. Deployment: The model is optimized for real-world use through compression, inference tuning and safety controls.

In summary,

LLMs are powerful language generators built by scaling data, compute and Transformer based architectures - but they operate on probability not genuine understanding.