

Core Concepts of Large Language Models (LLMs)

1. Introduction

Large Language Models (LLMs) such as GPT, Claude, and Gemini rely on a set of **core configuration concepts and parameters** that directly influence how responses are generated. Understanding these concepts is essential for **prompt engineering, application development, API integration, and responsible AI usage**.

This document explains the most important core concepts including **tokens, temperature, system roles, and other key parameters**, with examples and real-world software development use cases.

2. Tokens: The Fundamental Unit of LLM Processing

2.1 What Are Tokens?

A **token** is a chunk of text that an LLM processes. Tokens are **not the same as words**.

- A word may be split into multiple tokens
 - Punctuation and spaces can be tokens
 - Numbers and symbols are also tokenized
-

2.2 Why Tokens Matter

Tokens impact:

- Cost (API pricing is token-based)
 - Context window size
 - Performance and latency
 - Prompt and response length limits
-

2.3 Input Tokens vs Output Tokens

- **Input tokens:** Tokens sent in the prompt
- **Output tokens:** Tokens generated by the model
- **Total tokens:** Input + Output

2.4 Use Case: Token Awareness in Software Development

Scenario: A chatbot integrated into a customer support system must handle long conversations.

Solution:

- Limit conversation history
- Summarize older messages
- Control max output tokens

3. Temperature: Controlling Creativity and Randomness

3.1 What Is Temperature?

Temperature controls how **random or deterministic** the model's output is.

- 0.0 - 0.2: Very deterministic, factual
 - 0.3 - 0.5: Balanced
 - 0.6 - 0.8: Creative
 - 0.9 - 1.0: Highly creative, less predictable
-

3.3 Best Practices

- Code generation: 0.0 - 0.3
 - Chatbots: 0.4 - 0.6
 - Creative writing: 0.7 - 0.9
-

4. System, User, and Assistant Roles

4.1 Role-Based Prompting

LLMs support **role-based messages** to control behavior.

- **System:** Sets rules, tone, and behavior
 - **User:** Provides instructions or questions
 - **Assistant:** Model-generated response
-

4.2 System Role (Most Important)

The **system role** defines how the model should behave.

Example:

```
System: You are a strict JSON API. Return only valid JSON.  
User: Generate customer data.
```

5. Max Tokens: Limiting Response Length

Defines the **maximum number of tokens** the model can generate in a response.

Why It Is Important:

- Prevents long, unnecessary responses
 - Controls cost
 - Improves performance
-

6. Top-p (Nucleus Sampling)

Top-p limits token selection to the **most probable tokens whose cumulative probability $\leq p$** .

Best Practice: Use either temperature or top-p, not both aggressively.

7. Frequency and Presence Penalties

- **Frequency Penalty:** Reduces repeated words or phrases
 - **Presence Penalty:** Encourages introducing new topics
-

8. Practical Use Case: API-Based AI Assistant

Scenario: An organization builds an AI assistant for developers.

Configuration:

- System role: "You are a senior software architect"
 - Temperature: 0.3
 - Max tokens: 300
 - Top-p: 0.9
 - Frequency penalty: 0.2
-

9. Common Mistakes to Avoid

- Ignoring token limits
- Using high temperature for code generation
- Not defining system role
- Mixing explanation with structured output

- Overusing long prompts unnecessarily
-

10. Conclusion

Understanding core LLM concepts such as **tokens, temperature, roles, and generation parameters** is critical for building reliable, cost-effective, and production-grade AI applications. These parameters allow developers to control **accuracy, creativity, safety, and consistency**, making LLMs suitable for real-world software systems.