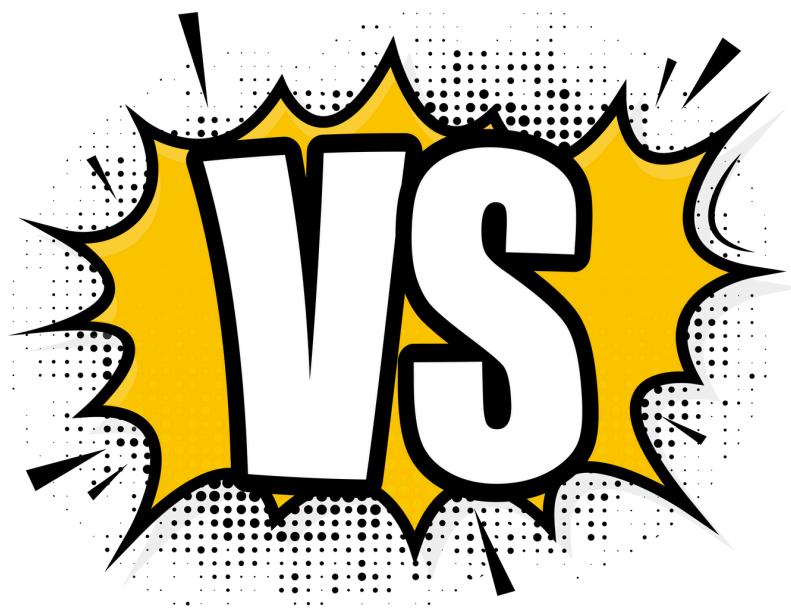


# Prompt Engineering



# Context Engineering



By SHAILESH SHAKYA @BEGINNERSBLOG.ORG



**Prompt Engineering** is the art and science of crafting precise instructions to get desired AI outputs. Think of it as "programming with words" for a probabilistic machine.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 1.1. Core Principles:

- Be Crystal Clear & Specific: Use action verbs. Define length, format, style, and audience. Avoid vague terms.
  - Bad: "Write about AI."
  - Good: "Explain the main differences between supervised and unsupervised learning in AI for a non-technical audience, in 3 bullet points."
- Provide Essential Context: Give relevant facts or background. Treat the AI like a brilliant, but uninformed, new team member.
  - Example: "Given that global temperatures have risen by 1 degree Celsius since the pre-industrial era, discuss the potential consequences for sea level rise."
- Show, Don't Just Tell (Examples): Provide input-output examples to define desired formats or styles. This is incredibly powerful.
  - Example: "Extract entities. Format: Company names: , People names: . Text: {text}"
- Iterate Relentlessly: Your first prompt won't be perfect. Test, observe, refine. It's a continuous loop.
- Frame Positively: Tell the AI "what to do," not "what not to do."
  - Bad: "Don't use passive voice."
  - Good: "Use active voice."
- Be Concise: Get to the point. Avoid unnecessary words, but don't sacrifice clarity.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 1.2. Fundamental Prompt Types:

Choose the right mode based on task complexity and how much guidance the AI needs.

### ● Zero-shot Prompting:

- **Concept:** Give the AI a task it hasn't seen before, with no examples. It relies purely on its vast pre-trained knowledge to infer a response.
- Simple, straightforward tasks where the AI's general knowledge is sufficient.
- **Example:**
  - Prompt: "Translate 'Hello, how are you?' to French."
  - AI Output: "Bonjour, comment allez-vous?"

### ● One-shot Prompting:

- Concept: Provide one example of the desired input-output format directly in the prompt. The AI learns from this single instance.
- When to Use: When you need a specific format or a simple classification, and one clear example is enough.
- Example (Sentiment Analysis):

#### ■ Prompt:

Review: "This movie was fantastic! I loved every minute."

Sentiment: Positive

Review: "The customer service was terrible."

Sentiment:

■ AI Output: "Negative"



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# ● Few-shot Prompting (In-context Learning):

- Concept: Include a small number of input-output examples in the prompt. This guides the AI on desired output structure, tone, and style, allowing it to generalize from limited data.
- When to Use: Complex content generation, technical domains, or when specific output requirements are crucial, and you don't have enough data for full fine-tuning.
- Example (JSON Extraction):

## ■ Prompt:

Text: "John Doe, 30, works at Acme Corp as a Software Engineer."

JSON: {"name": "John Doe", "age": 30, "company": "Acme Corp", "title": "Software Engineer"}

Text: "Jane Smith, 25, is a Marketing Specialist at Global Innovations."

JSON: {"name": "Jane Smith", "age": 25, "company": "Global Innovations", "title": "Marketing Specialist"}

Text: "Alice Brown, 35, is a Data Scientist at Tech Solutions Inc."

JSON:

■ AI Output: {"name": "Alice Brown", "age": 35, "company": "Tech Solutions Inc.", "title": "Data Scientist"}

- Caution: Be aware of potential biases (e.g., majority label bias, recency bias) if examples are skewed or ordered poorly.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# 1.3. Advanced Prompting Techniques:

These techniques push LLMs beyond basic responses, enhancing reasoning and control.

## ● Role-Playing / Persona Prompting:

- Concept: Assign a specific role or persona (e.g., "senior legal analyst," "fantasy author") to the AI. This influences its style, tone, focus, and depth of information.
- Why it Works: The AI taps into its vast training data to mimic communication patterns and domain knowledge associated with that role.
- Example:
  - Prompt: "You are a university professor specializing in machine learning. Explain the concept of neural networks to first-year students."
  - AI Output: (Will explain in an academic, structured, yet accessible way, like a professor.)



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# ● Chain-of-Thought (CoT):

- Concept: Guide the AI to show its intermediate reasoning steps before giving the final answer.
- Improves logical reasoning, mathematical computations, and transparency, especially for complex tasks. It makes the AI "think aloud."

## ○ Example:

- Prompt: "A store has 12 apples. It sells 5, then buys 3 more. How many apples are left? Let's think step by step."
- AI Output: "First, the store had 12 apples. It sold 5, so  $12 - 5 = 7$  apples. Then, it bought 3 more, so  $7 + 3 = 10$  apples. The answer is 10."



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## ● Tree-of-Thought (ToT):

- Concept: An advanced CoT. The AI explores multiple distinct reasoning paths, self-evaluates choices, and can backtrack.
- Enables deliberate decision-making and planning for tasks requiring non-trivial search or creativity (e.g., complex puzzles, creative writing).

## ● Self-Consistency:

- Concept: The AI samples a diverse set of reasoning paths for a problem and then selects the most frequent answer among them (majority vote).
- Improves reliability and robustness by aggregating multiple "guesses," leveraging the AI's ability to judge its own outputs.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## ● Least-to-Most (LtM):

- Concept: Breaks down complex problems into simpler subproblems and solves them sequentially, explicitly using the output of each subproblem as input for the next.
- Addresses AI's "working memory" limitations by managing intermediate states, allowing it to generalize to much longer reasoning chains.

## ● Generated Knowledge:

- Concept: The AI first generates relevant knowledge itself about the task, then uses that generated knowledge as additional context to formulate the final answer.

Enhances contextual understanding and accuracy by self-augmenting its knowledge, reducing misinterpretation.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# 1.4. Structuring & Automating Prompts

## ● Prompt Chaining:

- **Concept:** Link multiple prompts in a logical sequence, where the output of one becomes the input for the next.
- Breaks down complex tasks into manageable subtasks, improving control, reliability, and consistency in workflows.
- **Example:** Prompt 1: "Summarize this document."  
-> Output 1. Prompt 2: "Rephrase the summary from Output 1 for a marketing email." -> Output 2.

## ● Prompt Templating:

- **Concept:** Standardize instructions and input formats using placeholders (variables).
- Ensures consistency, clarity, efficiency, and reusability across tasks or teams. Think of it like code libraries for prompts.

**Tool:** LangChain's PromptTemplate allows easy creation and management of templates.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# 1.5. Practical Applications

- **Summarizing:** Condense long texts.
- Types: Extractive (selects original sentences) vs. Abstractive (generates new sentences).
- **Tip:** Be specific about length and focus.
  
- **Rewriting:** Adjust style, tone, or length.
- **Action:** "Rewrite this paragraph in a formal tone."
- **Tip:** Provide your own writing samples for style mimicry.
  
- **Generating Agent Instructions:** Define the AI's role, personality, planning, memory, and tool usage for complex tasks.
- **Example:** "You are a helpful code assistant that can teach a junior developer how to code. Your language of choice is Python. Don't explain the code, just generate the code block itself."
  
- **Code Generation:** Generate code from comments, complete functions, or create database queries.
- **Tip:** Provide clear requirements and examples.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



When using LLM APIs (OpenAI, Google AI Studio, Anthropic, Hugging Face), these parameters control output behavior.

- **temperature:** Controls creativity/randomness. Lower (e.g., 0.2) for factual, higher (e.g., 0.8) for creative.
- **top\_p:** Similar to temperature, controls determinism. Lower for exact answers, higher for diverse responses.
- **max\_tokens / max\_length:** Sets the maximum length of the AI's response, controlling cost and verbosity.
- **stop sequences:** Strings that tell the AI to stop generating (e.g., "11." to limit a list to 10 items).

**frequency penalty / presence penalty:** Reduce word/phrase repetition



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# Part 2: Context Engineering

Context Engineering is the strategic design of the entire information environment an AI model perceives before generating a response. It's about building a "mental world" for the AI, not just giving it a single instruction.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 2.1. The Big Picture: LLM as CPU, Context as RAM

- **Analogy:** LLM is like a CPU of an AI system, and its context window is the RAM.
- **Prompt Engineering:** What you tell the CPU to do right now (a single instruction).
- **Context Engineering:** How you manage all the data in the RAM (context window) so the CPU has everything it needs to perform any task, consistently, over time.
- A brilliant prompt is useless if the AI's "RAM" is filled with irrelevant data or lacks crucial information. Context Engineering ensures the AI always has the right "knowledge" at its fingertips.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 2.2. Context Engineering vs. Prompt Engineering:

A Clear Distinction Prompt Engineering is a subset of Context Engineering.

Aspect	Prompt Engineering	Context Engineering
Focus	Crafting the perfect single instruction ("What to say").	Designing the entire information environment ("What the AI knows, and why it cares").
Scope	Narrow: Immediate input prompt.	Broad: Memory, history, tools, system prompts, external data.
Mindset	Creative writing, linguistic optimization.	Systems design, information architecture.
Scalability	Hard to scale; needs manual tweaks for variations.	Built for scale; handles complex, long-running workflows.
Debugging	Rewording prompts, guessing.	Inspecting full context window, memory, data flow.
Effort Type	Linguistic creativity.	Architectural thinking, data management.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 2.3. Context Window Management Strategies

LLMs have a finite "working memory" (context window) measured in tokens. Exceeding it degrades performance.

- **The Problem:** Large context windows can introduce noise, increase cost, and hide crucial "needles in a haystack."

### ● Strategies:

- **Chunking & Summarization:** Break large texts into smaller, manageable segments (chunks). Summarize accumulated information to preserve details and remove redundancy.
- **Sliding Window Technique:** Process text in overlapping segments, focusing attention on a fixed-size window. Reduces computational complexity while maintaining context continuity.
- **Filtering Irrelevant Information:** Actively remove unnecessary details or "noise" before they reach the model.
- **Dynamic Contextualization:** Adjust context window size and content based on the specific query.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 2.4. Memory Mechanisms: Short-term vs. Long-term

To maintain coherent interactions, AI systems need different types of memory.

### ● Short-Term Memory (STM):

- Concept: The immediate context window of the LLM. It's ephemeral; information disappears once the response is generated.
- Purpose: Holds information relevant to the current API request or conversation turn. Crucial for real-time tasks.
- Limitation: Limited capacity (around 5-9 pieces of info), leading to context loss if not managed.

### ● Long-Term Memory (LTM):

- Concept: External systems designed to store data beyond a single chat session, allowing the AI to "remember" details, user preferences, or vast document repositories over extended periods.
- Mechanisms: Often uses external vector databases and graph databases to store information in numerical formats that capture relationships.
- Benefits: Enables context-rich responses, significantly reduces hallucinations (via RAG), and facilitates efficient data handling. It's key for "lifelong learning" in AI.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# Part 3: Building Robust AI Systems with Context Engineering (RAG)

Retrieval-Augmented Generation (RAG) is a cornerstone for building knowledgeable and reliable LLM applications, overcoming their static knowledge and tendency to "hallucinate."



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 3.1. RAG Workflow:

RAG dynamically fetches external information to enhance the LLM's output with up-to-date, verifiable, or domain-specific facts.

### 1. Indexing (Offline):

- Action: Load raw data (documents, chat history), preprocess (chunk into smaller segments), convert to high-dimensional vector embeddings using an embedding model, and store in a vector database.
- Why: Creates a searchable knowledge base.

### 2. Retrieval (Runtime):

- Action: User query is converted into a vector embedding. A similarity search finds the most semantically relevant documents/chunks from the vector database.
- Why: Finds the "needle in the haystack" – the most relevant information for the query.

### 3. Generation:

- Action: The retrieved relevant documents are combined with the original user prompt. This augmented input is fed to the LLM.
- Why: Ensures the AI's response is accurate, context-aware, and grounded in real data, mitigating hallucinations.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 3.2. Key Components: Embeddings & Vector Databases

### ● Embeddings:

- Concept: Numerical representations of text (words, sentences, documents) in a high-dimensional space. They capture semantic meaning and relationships.
- Why: Allow for "semantic search" – finding information based on meaning, not just keywords.

### ● Vector Databases:

- Concept: Specialized databases designed to store embeddings and facilitate rapid similarity searches.
- Examples: Pinecone, Weaviate, Milvus, Qdrant. Traditional databases like Redis, pgvector, Elasticsearch also offer vector support.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 3.3. Indexing & Retrieval Strategies

Crucial for efficient and accurate information retrieval in RAG.

- **Chunking Strategies:** How large documents are split (e.g., fixed-length, sentence-level, recursive, content-aware).
- **Indexing Mechanisms:** How vectors are organized for fast search (e.g., HNSW, IVF, PQ, LSH).
- **Retrieval Methods:**
  - **Hybrid Search:** Combines keyword and semantic search.
  - **Recursive Retrieval:** Starts with small chunks, then retrieves larger ones for richer context.
  - **Sub-Queries:** Breaks complex user queries into smaller questions for different data sources.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 3.4. Integrating with LLM Frameworks

Frameworks simplify building RAG pipelines, abstracting complexity.

- **LangChain:** Provides PromptTemplate options and facilitates prompt chaining.
- **LlamaIndex:** Powerful open-source framework specifically for RAG, simplifying chunking, embedding, indexing, and querying. Offers "recipes" for optimization.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# Part 4: Evaluating, Debugging & Optimizing LLM Outputs

Essential for ensuring reliable and high-quality AI systems.

## 4.1. LLM Evaluation Metrics

Assess output quality across correctness, relevance, and efficiency.

- **Quantitative:** Accuracy, Precision, Recall, F1-Score (classification); BLEU, ROUGE (text generation); Perplexity.
- **LLM-as-a-Judge:** Use one LLM to evaluate another's output based on natural language rubrics (e.g., G-Eval). Often more accurate for semantic nuance.
- **RAG-Specific:** Retrieval Accuracy, Relevance Score, Response Coherence, Factual Consistency, Robustness to Noise, Negative Rejection, Information Integration.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 4.2. Common Failure Modes:

Understanding these helps you build resilient AI.

- **Prompt Engineering Failures:** Vague prompts, prompt overload, ambiguous language, missing context, lack of iteration, blind trust in AI output, forgetting role/audience, biases.
- **RAG System Failures:** Extraction errors (AI misreads retrieved info), context size limitations (info too long), inexhaustive computation (AI responds too early), missing content (info not in database), missed top-ranked (relevant info not retrieved), wrong format/specificity, answer not extracted (AI fails to find answer in provided context).



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# 4.3. Debugging & Iterative Prompt Design

A systematic process for optimizing performance.

- 1. Spot Problems:** Identify hallucinations, unclear instructions, token limit issues.
- 2. Test Each Part:** Verify input data, context clarity. Start simple, add complexity.
- 3. Use Tools:** Visualize attention maps, analyze token probabilities, use tracking software (LangSmith, PromptLayer).
- 4. Iterate & Improve:** Track prompt versions, test variations, incorporate user feedback.
- 5. Final Testing:** Measure performance with metrics, test real-world scenarios, human review.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 4.4. A/B Testing for Prompt Optimization

- **Concept:** Systematically compare different prompt versions to see which performs better.
- Essential for continuous evaluation, enhancing performance, reducing bias, accelerating innovation, and data-driven decision-making.
- **Action:** Integrate A/B testing into continuous deployment workflows.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# Part 5: Ethical Considerations & Responsible AI Practices

Crucial for safe, fair, and transparent LLM deployment.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# 5.1. Bias in LLM Outputs & Mitigation

- **Types:** Demographic, language/cultural, temporal, stereotypical association. Often reflect training data flaws.
- **Mitigation:**
  - Prompt Design: Use clear, neutral prompts; include diverse examples; explicitly instruct AI to avoid stereotypes.
  - Advanced: Reinforcement Learning from Human Feedback (RLHF), adversarial prompting, context isolation.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 5.2. Hallucinations & Ethical Implications

- **Definition:** AI generates factually incorrect, inconsistent, or made-up text, despite sounding plausible. Not human error; stems from statistical pattern prediction.
- **Causes:** Errors in encoding/decoding, data outliers, model prioritizing fluency over factual accuracy.
- **Implications:** Undermines reliability/trust, spreads misinformation, amplifies bias, privacy risks, reputational/financial damage, safety concerns (e.g., in healthcare).
- **Mitigation:** RAG (grounding in verifiable facts), robust evaluation, continuous monitoring, guardrails (detect and filter unsafe outputs).



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



## 5.3. Responsible AI Practices & Governance

Ensuring safe, secure, fair, and ethical use of LLMs.

### ● Key Aspects:

- **Fairness:** Objective systems, no discrimination.
- Accountability & Human Oversight: Clear roles, human intervention for high-risk applications.
- **Privacy & Data Security:** Encryption, anonymization, access controls. Compliance with laws (HIPAA, CCPA).
- **Security:** Protect from unauthorized access/abuse (adversarial attacks). Use input firewalls, output validators.
- **Transparency:** Be clear about AI capabilities, limitations, decision-making.
- **Continuous Monitoring:** Regular assessment for fairness, accuracy, compliance.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# Part 6: Future Trends

The AI field evolves rapidly. Stay ahead.

## ● Anticipated LLM Advancements:

- Enhanced contextual understanding (e.g., GPT-4o).
- Multimodal AI (processing text, images, audio, video).
- Adaptive prompting (AI adjusts to user style).
- Increased reliability and efficiency (less hallucinations, faster, cheaper).

## ● Emerging Prompt/Context Trends:

- **AI for Prompt Creation:** Generative AI creating and optimizing prompts.
- **"Mega-Prompts":** Using much larger context windows (e.g., Gemini 1.5 Pro's 1M tokens).
- **Automated Prompt Optimization:** ML algorithms (RL, Bayesian optimization) to find best prompts.
- **Context Engineering as Primary Skill:** Shift from single prompts to managing the entire AI information ecosystem.
- **Human-AI Collaboration:** AI as a brainstorming partner, refining drafts.
- **No-code Platforms:** Empowering non-technical users to create prompts.
- **Ethical Prompting:** Continued focus on fairness, transparency, bias mitigation.



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG



# ● Learning:

- **Stay Updated:** Read academic papers (NeurIPS, ACL, EMNLP, arXiv).
- **Engage Communities:** Participate in forums (e.g., Reddit's r/PromptEngineering).
- **Online Courses:** Enroll in specialized courses (Coursera, Udemy, edX from Google, IBM, Vanderbilt, DeepLearning.AI).
- **Hands-on:** Experiment with playgrounds and APIs (OpenAI, Google AI Studio).
- **Frameworks:** Master LangChain, LlamalIndex.
- **Engineering Discipline:** Apply version control for prompts, A/B testing, robust evaluation.  
Problem Formulation: Focus on clearly defining the problem's scope, not just prompt phrasing



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG





I hope you have found this information helpful

Join **OpenAI Learning** to get more educational stuff Similar to this you finished reading  



Telegram: [OpenAI Learning](#)



WhatsApp: [OpenAI Learning](#)

**Thank You!**



By SHAILESH SHAKYA



POWERED BY:  
BEGINNERSBLOG.ORG

Swipe to  
Next Slide 





Created by Shailesh Shakya  
**@BEGINNERSBLOG.ORG**

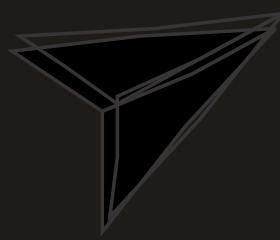
Did you find this post helpful?  
Please...



LIKE



COMMENT



REPOST



SAVE