# DAY 6

# Hands-On with Gemini Playground AI: Prompt Engineering & Advanced Features

## What is Gemini Playground AI?

Gemini Playground (also known as Google AI Studio) is a tool where you can interact with Gemini models in real-time, test prompt variations, and understand how AI generates responses. Unlike standard chatbots, Gemini gives you full control over the environment - you can define how the model thinks, responds, formats answers, and even calls external functions.

Think of it as an AI sandbox — a place to play, test, learn, and build.

## Flash vs Pro: Gemini 2.5 Model Variants

Gemini now offers multiple versions:

- ❖ Gemini 2.5 Pro – Full-power model for rich responses and deep reasoning.
- ❖ Gemini 2.5 Flash – A faster and lighter version designed for low-latency outputs, especially in real-time applications.

I used Gemini 2.5 Flash for most tests today. It's incredibly responsive and still quite accurate for summaries, explanations, and role-based prompting.

## System vs User Prompts

Prompts in Gemini Playground are divided into two categories:

**System Prompt**

❖ Sets the *role* or *personality* of the model.

❖ Invisible to the user but strongly affects tone, detail, and attitude.

❖ Example:

"You are a calm, empathetic teacher who explains clearly.

**User Prompt**

❖ The visible input from the user.

❖ Tell the AI what the *task* is.

❖ Example:

"Explain Newton's Laws to a 10-year-old with examples."

Combining both helps simulate different characters, tones, or industries — from doctors to poets to sarcastic assistants.

## Temperature & Top-p – Controlling Randomness

**Temperature:**

| Temp Range | Behavior |
|---|---|
| 0.0 - 0.3 | Very precise and fact-based; little variation. |
| 0.4 - 0.7 | Balanced; good for informative or explanatory content. |

| 0.8 -1.3 | More imaginative, expressive, creative language. |
|---|---|
| 1.4 - 2.0 | Highly creative, sometimes random or surprising responses. |

❖ Use lower temperatures for answers you want to be stable and fact-based.

❖ Use higher temperatures if you want the AI to think out of the box, give multiple

perspectives, or generate creative writing.

**Top-p:**

Nucleus sampling — limits output to the top % of likely tokens.

E.g. Top-p = 0.95 → Only the most probable words are considered.

Used together, these give precise control over how "safe" or "creative" the AI should be.

**Test Scenarios**

| Test | System Prompt | User Prompt | Temperature | Output Summary |
|---|---|---|---|---|
| 1. | None | Tell me about AI in agriculture | 0.2 | Accurate, academic-style explanation. No creative flair. |
| 2. | You are a calm and professional teacher. | Explain the water cycle to a 5th grader. | 0.3 | Gave a clear, step-by-step breakdown. Very factual. |

| 3. | You are a motivational speaker. | Give me tips to stop procrastinating. | 0.7 | Delivered energetic, action-oriented advice. Easy to follow. |
| --- | --- | --- | --- | --- |
| 4. | You are a calm and professional teacher. | Explain the water cycle to a 5th grader. | 1.2 | Added a playful tone, used a storytelling approach ("Imagine you're a water droplet…"). |
| 5. | You are a sarcastic assistant. | Why is homework amazing? | 1.8 | Humorous, ironic tone with creative exaggeration. Nailed the sarcasm. |

## Observations

❖ System prompts strongly affect tone and style—Gemini adapts well to different personas (e.g., poet, coach, teacher).

❖ Temperature at 1.2+ adds expressive flair but may skip detailed logic or structure.

❖ Without system prompts the tone is neutral and informative but lacks flavor.

❖ Creative writing prompts perform best at **1.3 – 1.8** temperature range.

❖ Gemini responds faster at low temps but more richly at higher ones.

❖ Lower temps = faster responses; higher temps = richer language.

## Tokens & Output Length

Gemini works with tokens, just like OpenAI's models. A token is a chunk of text (usually ~4 characters).

**Example:** "AI is fun." → ["AI", " is", " fun", "."] → 4 tokens

**Token Limit (Gemini 2.5):**

- ❖ Up to 1,048,576 tokens

- ❖ My test prompt: 60–100 tokens

- ❖ Max output length setting: 65,536 tokens

**Why it matters:**

- ❖ Tokens affect cost, latency, and memory

- ❖ Long prompts + big output = more tokens used

- ❖ Helps when building long chat flows or structured apps

## Thinking Mode & Thinking Budget

This is a Gemini-exclusive experimental feature.When enabled, the AI:

- ❖ Pauses slightly before responding

- ❖ Internally "thinks" through the task

- ❖ Outputs a more reasoned, step-by-step answer

You can adjust the Thinking Budget (how long the AI thinks). Great for:

- ❖ Chain-of-thought reasoning

- ❖ Math problems

❖ Logical planning

**Example:** Prompt: "Solve this step-by-step: A train leaves at 40km/h…"

Gemini (with Thinking Mode ON) returned:

❖ Each calculation step

❖ Realistic timeline

❖ Final answer with reasoning

## Structured Output (JSON / Markdown / Tables)

This feature is extremely useful for developers, data analysts, and automation workflows.

You can prompt Gemini to:

❖ Format responses in JSON

❖ Create Markdown tables

❖ Output clean HTML

❖ Or follow custom schemas

**Example:**

"List 5 programming languages in a JSON array with fields: name, use case, difficulty."

**Gemini output:**

```
[
  { "name": "Python", "use_case": "AI, scripting", "difficulty": "Easy" }, ...
```

]

Super helpful for:

- ❖ Resume generators

- ❖ Report formatting

- ❖ API data generation

- ❖ Structured extraction from messy input

## Function Calling (Beta)

Gemini supports function calling, similar to OpenAI tools. You can:

- ❖ Simulate plugin-like behavior

- ❖ Let Gemini call predefined backend functions

- ❖ Control the flow of logic between AI ↔ system

Still experimental, but very promising for:

- ❖ Tool use

- ❖ Calendar lookups

- ❖ Live data integrations

## Code Execution

Gemini also supports code execution in the Playground (when enabled).This means you can:

- ❖ Run Python snippets directly

- ❖ Test code logic or math functions

❖ Automate reasoning-based prompts

**Example:**

Prompt: "Write and run a Python function that checks if a number is prime."

**Result:**

❖ Gemini generated the code

❖ Ran the logic and returned output instantly

## Google Grounding & URL Context

This is one of the coolest features:

❖ You can ground Gemini's answers using live Google Search

❖ Or provide a URL context (paste a link to a website)

**Benefits:**

❖ More accurate, up-to-date results

❖ Perfect for research, product comparisons, fact-checking

**Example:**"Using URL context: https://en.wikipedia.org/wiki/Photosynthesis — summarize in 5 bullet points."

Gemini scanned the page and responded with a real-time summary. Amazing!

## Additional Tools & Settings

| Setting/Feature | Description |
|---|---|
| Safety Settings | Control response filtering based on content type |
| Stop Sequences | Tells Gemini where to stop generating |
| Memory & Chat History | Saves previous turns in conversations(when supported) |
| Advanced Parameters | CFG Scale,Response Timeout, Role Templates |

## What Is MCP Server?

While not directly visible in Gemini UI, MCP (Message Control Protocol / Management Control Platform) may refer to:

- ❖ Backend servers that manage session flow, user control, access limits

- ❖ Enterprise environments using Gemini in a secure, controlled space

Think of it as the infrastructure layer — not the model itself.

## Final Takeaways

- ❖ Gemini Playground is a next-level platform for exploring prompt engineering.

- ❖ From role prompting to structured JSON output — it gives full control.

- ❖ Thinking Mode, code execution, and function calling make it ideal for serious tasks.

- ❖ Features like Google grounding and URL context make Gemini smarter than most static models.

- ❖ With such depth and flexibility, Gemini can be used in education, app development, research, writing, and AI-powered automation.