# Google's Gemini 2.5 Cuts AI Costs by 600% with 'Thinking Budgets'

## Description

Google has recently unveiled Gemini 2.5 Flash, a significant enhancement to its AI platform, providing businesses and developers remarkable control over the model's cognitive capabilities. This new version, made available today for preview via Google AI Studio and Vertex AI, aims to boost reasoning skills while remaining competitively priced in a saturated market.

A key feature of this model is the "thinking budget," which permits developers to set the computational resources allocated for complex reasoning before generating a response. This innovation addresses a crucial challenge in AI: the trade-off between sophisticated reasoning and increased latency and costs.

"Understanding the significance of cost and latency for developers, we are committed to providing the flexibility to adjust the model's cognitive processes," stated Tulsee Doshi, Google DeepMind's Product Director for Gemini Models, in an exclusive interview.

This flexibility highlights Google's pragmatic approach to AI, allowing businesses to enable or disable reasoning based on specific needs. This advancement, which they term their "first fully hybrid reasoning model," gives users unprecedented control over AI deployment.

## Vocabulary List:

1. **Enhancement** /ɪnˈhænsmənt/ (noun): An improvement or increase in quality value or extent.
2. **Cognitive** /ˈkɒɡnɪtɪv/ (adjective): Relating to the processes of thought and understanding.
3. **Reasoning** /ˈriːzənɪŋ/ (noun): The action of thinking about something in a logical sensible way.
4. **Latency** /ˈleɪtənsi/ (noun): The delay before a transfer of data begins following an instruction for its transfer.
5. **Pragmatic** /præɡˈmætɪk/ (adjective): Dealing with things sensibly and realistically in a way that is based on practical rather than theoretical considerations.
6. **Deployment** /dɪˈplɔɪmənt/ (noun): The action of bringing resources into effective action.

## Comprehension Questions

**Multiple Choice**

1. What is the name of the new AI platform enhancement recently unveiled by Google?

   Option: A. Gemini 2.0 Flash
   Option: B. Gemini 2.5 Burst
   Option: C. Gemini 2.5 Flash
   Option: D. Gemini 3.0 Spark

2. Where can developers preview the new Gemini 2.5 Flash version?

   Option: A. Google Drive
   Option: B. Google AI Studio
   Option: C. Google Maps
   Option: D. Google Play Store

3. What is the key feature of Gemini 2.5 Flash that allows developers to allocate computational resources?

   Option: A. Computing Budget
   Option: B. Cognitive Reserve
   Option: C. Thinking Budget
   Option: D. Intelligent Quota

4. Who is quoted as stating the significance of cost and latency for developers regarding Gemini 2.5 Flash?

   Option: A. Sundar Pichai
   Option: B. Larry Page
   Option: C. Tulsee Doshi
   Option: D. Eric Schmidt

5. What is one of the benefits of Gemini 2.5 Flash highlighted in the text?

   Option: A. Increased Latency
   Option: B. Reduced Reasoning Skills
   Option: C. Greater Control Over AI Models
   Option: D. High Costs

6. What is the term used by Google to describe Gemini 2.5 Flash as their innovation?

   Option: A. Cognitive Revolution
   Option: B. Hybrid Evolution
   Option: C. Fully Automated Model
   Option: D. First Fully Hybrid Reasoning Model

**True-False**

7. Gemini 2.5 Flash was designed to reduce reasoning skills in AI models.

8. The trade-off addressed by Gemini 2.5 Flash involves increased latency and costs in AI.

9. According to the text, businesses can now disable reasoning in AI models based on their specific needs.

10. Tulsee Doshi is the Product Director for Google DeepMind.

11. Google AI Studio and Vertex AI are the only platforms where Gemini 2.5 Flash can be previewed.

12. Gemini 2.5 Flash is the first fully hybrid reasoning model offered by Google.

**Gap-Fill**

14. The flexibility offered by Google in adjusting the model's cognitive processes highlights their pragmatic

approach to _____ .

15. Businesses now have the ability to enable or disable reasoning in AI models based on their specific

_____ .

16. The thinking budget feature of Gemini 2.5 Flash addresses the crucial challenge in AI: the trade-off

between sophisticated reasoning and increased _____ .

17. Gemini 2.5 Flash gives users unprecedented control over AI deployment as it is the

_____ offered by Google.

18. Developers can preview Gemini 2.5 Flash via Google AI Studio and _____ .

# Answer

**Multiple Choice:** 1. C. Gemini 2.5 Flash 2. B. Google AI Studio 3. C. Thinking Budget 4. C. Tulsee Doshi
5. C. Greater Control Over AI Models 6. D. First Fully Hybrid Reasoning Model
**True-False:** 7. False 8. True 9. True 10. False 11. False 12. True
**Gap-Fill:** 14. AI 15. needs 16. latency 17. first fully hybrid reasoning model 18. Vertex AI

**CATEGORY**

1. Sci/Tech - LEVEL4

**Date Created**
2025/04/18
**Author**
aimeeyoung99