

## Most frequently used LLMs June 2025

Model	Organization	Specialization	Training Data	Max Context	Cost for Pro Plan
Microsoft 365 Copilot, Copilot in GitHub	Microsoft/GitHub	Code-centric AI assistant in IDEs, Generalist conversational	N/A	Uses GPT-4.5 Input- 128,000 tokens Output- 16,000 tokens	\$30/user/month Microsoft 365 Copilot \$20/user/month Copilot in GitHub,
Llama 4 (Open Source)	Meta	Open-source multimodal/text-code model	15 trillion tokens	Input- 128k tokens (405B) Output- 2,048 tokens	Token based-Range \$0.30-5.33/1M input tokens
OpenAI (GPT-4.5) (o3)	OpenAI	Generalist conversational and coding model	Proprietary (estimated trillions)	Input- 128,000 tokens Output- 16,000 tokens	\$20/month
Claude	Anthropic	Ethical, long-context conversation and reasoning	Proprietary	Input- 128,000 tokens Output- 64,000 tokens	\$20/month
Gemini	Google	Multimodal, efficient transformer TPU-optimized	Multimodal web/code datasets	Input- 1,000,000 tokens Output- 8,000 tokens	\$19.99/month
Granite	IBM	Open-source multimodal/text-code model	Unknown	Input- 128,000 tokens Output- 16,384 tokens	\$140 per month for 100 users
Qwen	Alibaba (Chinese based)	Not enough public info yet	Unknown	Unknown	\$24.99 - \$39.99/month
Grok	xAI	Real-time web-connected chatbot	Trained on web/X (Twitter) data	Input-131K tokens Output- 4,096 tokens	\$30-40/month
DeepSeek	High-Flyer (Chinese based)	Has specialized coder, math, and vision models	Unknown	Unknown	Token based-Input \$0.55-0.14/1M input tokens Output <b>\$2.19</b> /1M output tokens

One Token= Approximately 4 English language characters

Input=User Prompt

Output=Response to prompt