



2025年11月最新生成AIモデル徹底比較レポート： MiniMax-M2、Kimi K2 Thinking、ERNIE-4.5-VL- 28B-A3B-Thinking、GPT-5.1の性能分析

本レポートでは、2025年10月下旬から11月中旬にかけて立て続けにリリースされた4つの最先端生成AIモデル——**MiniMax-M2**（10月26日）、**Kimi K2 Thinking**（11月6日）、**ERNIE-4.5-VL-28B-A3B-Thinking**（11月11日）、**GPT-5.1**（11月12日）——について、公式技術レポート、第三者ベンチマーク、および業界専門家の評価を総合的に分析し、各モデルの性能特性、実用的仕様、そして2025年11月時点における生成AI技術の最前線を明らかにする。[\[1\]](#) [\[2\]](#) [\[3\]](#) [\[4\]](#) [\[5\]](#)

これら4モデルのリリースは、生成AI業界における重要な転換点を示している。特に中国勢（**MiniMax**、**Moonshot AI**、**Baidu**）が相次いでオープンソースの高性能モデルを投入し、従来OpenAIやAnthropicが独占していた技術的優位性に挑戦する構図が鮮明となった。同時に、OpenAIもGPT-5.1で会話性とユーザビリティの大幅改善を図り、市場の競争激化に対応している。[\[2\]](#) [\[3\]](#) [\[6\]](#) [\[7\]](#) [\[5\]](#) [\[8\]](#) [\[9\]](#)

1. MiniMax-M2：コスト効率と推論速度の革新（2025年10月26日リリース）

1.1 技術アーキテクチャと基本仕様

MiniMax-M2は、中国のAIスタートアップMiniMaxが開発した**Mixture-of-Experts (MoE)**アーキテクチャを採用する大規模言語モデルである。総パラメータ数は**230億（230B）に達するが、推論時に活性化されるのはわずか10億パラメータ（10B）**のみという極めて効率的な設計となっている。この設計により、大規模モデルの性能を維持しながら、推論速度と運用コストを劇的に削減することに成功している。[\[10\]](#) [\[11\]](#) [\[12\]](#) [\[1\]](#) [\[2\]](#)

モデルは**Apache 2.0**ライセンスの下でオープンソース化されており、Hugging Face上で重みが公開され、vLLMやSGLangなどの推論フレームワークでのデプロイが可能である。コンテキスト윈드ウは**204,800トークン（約204K）をサポートし、出力容量は131,072トークン（約131K）**に達する。この広大なコンテキスト윈드ウにより、複数ファイルにまたがるコードベース全体や長文ドキュメントの処理が可能となっている。[\[11\]](#) [\[13\]](#) [\[14\]](#) [\[2\]](#) [\[10\]](#)

技術的な特徴として、MiniMax-M2は独自の「Lightning Attention」メカニズムを実装しており、高速かつ低メモリでの推論を実現している。MoEアーキテクチャでは、タスクの内容に応じて必要な「専門家（エキスパート）」のみを選択的に活性化させることで、計算効率を最大化している。実際の運用環境では、NVIDIA H100 GPUをわずか4枚使用するだけで高性能な推論が可能であり、中規模組織でも導入可能なハードウェア要件となっている。[\[15\]](#) [\[1\]](#) [\[2\]](#)

1.2 ベンチマーク性能：オープンソース最高峰の実力

MiniMax-M2は複数の権威あるベンチマークテストで驚異的な成績を記録している。最も注目すべきは、**Artificial AnalysisのIntelligence Index v3.0**において**61ポイント**を獲得し、**オープンソースモデルとして史上最高スコア**を達成したことである。この成績により、グローバル全体でも第5位にランクインし、上位にはGPT-5、Grok 4、Claude Sonnet 4.5といった最先端のプロプライエタリモデルのみが位置している。[\[12\]](#) [\[1\]](#) [\[2\]](#)

コーディング能力では**SWE-bench Verified**で**69.4%**のスコアを記録し、実務的なプログラミング課題において高い正答率を示した。これはGPT-4oの約**65%**を大きく上回る成績である。さらに、ツール利用能力を測定する**T²-Bench**では****77.2%****というスコアで、業界平均を20%以上上回る卓越した性能を発揮している。[\[16\]](#) [\[1\]](#) [\[2\]](#)

エージェント型ワークフローに関連するベンチマークでも強みを見せており、**Terminal-Bench**では**46.3%**、**BrowseComp**では****44.0%****のスコアを記録した。これらの評価は、コマンドライン操作、ブラウジング、情報検索といった実務的なエージェントタスクにおける能力を示している。[\[17\]](#) [\[2\]](#)

独立検証テストでは**95%の正答率**を達成し、GPT-4oの約90%を上回る精度を示している。特にマルチファイル編集、コンパイル-実行-修正のループ、テスト駆動開発などのコーディングワークフローにおいて優れた性能を発揮することが確認されている。[\[14\]](#) [\[1\]](#) [\[10\]](#)

1.3 API価格と推論速度：コスト革命の実現

MiniMax-M2の最大の革新は、**圧倒的なコスト効率**にある。公式API価格は、入力トークンが**100万トークンあたり0.30ドル**、出力トークンが**100万トークンあたり1.20ドル**という極めて低価格に設定されている。この価格設定は、**Claude Sonnet 4.5のわずか約8%のコスト**に相当し、業界に衝撃を与えている。[\[1\]](#) [\[2\]](#) [\[16\]](#) [\[11\]](#)

推論速度においても画期的な性能を達成しており、**毎秒約100トークン**の生成速度を実現している。これはClaude Sonnet 4.5やGPT-4oと比較して**約2倍の速度**であり、実際の開発ワークフローにおいて「書く→編集→再生成→改良」のサイクルを大幅に高速化できる。初回トークンの到達時間も極めて短く、画像からJSONへの変換タスクでGPT-4oの約0.7倍の時間で完了することが確認されている。[\[16\]](#) [\[1\]](#)

コスト効率の実例として、開発者コミュニティでは「100万トークンの処理がわずか1.50ドルで可能」という報告があり、大規模な自動化パイプラインやバッチ処理において劇的なコスト削減を実現している。ある開発チームの報告では、MiniMax-M2への切り替えにより、API利用コストを前年比で**約92%削減**できたという事例も存在する。[\[17\]](#) [\[16\]](#)

1.4 実用的な強みと制約

MiniMax-M2の主な強みは、**エージェント型ワークフローとコーディング支援**に特化した設計にある。実際のソフトウェア開発現場において、コード生成、デバッグ、マルチファイル編集、APIツール連携などのタスクで高い実用性を示している。特に、コスト制約のあるスタートアップや中小企業にとって、**フラッグシップモデル並みの性能を低成本で利用できる点**は大きな魅力となっている。[\[2\]](#) [\[10\]](#) [\[14\]](#) [\[16\]](#)

一方で、いくつかの制約も報告されている。多言語タスクにおいては、**中国語タスクで80%の精度**に対し、Claude 3.5 Sonnetの85%を下回る性能となっている。また、創造的なナラティブ生成では

「言語ブレンドが多発する」という指摘があり、純粋な文学的創作には向かない可能性がある。さらに、エコシステムがまだ発展途上であり、GPT-4oやClaudeと比較してサードパーティツールやプラグインの統合が限定的である。[\[18\]](#) [\[16\]](#) [\[17\]](#)

ハードウェア要件としては、FP8精度での推論にNVIDIA H100 GPU 4枚が推奨されており、ローカル環境でのデプロイには相応のインフラ投資が必要となる。しかし、APIサービスを利用すれば、この制約を気にすることなく低コストで利用できるため、多くのユーザーにとってはホスティッドAPIが現実的な選択肢となっている。[\[19\]](#) [\[10\]](#) [\[11\]](#) [\[2\]](#)

2. Kimi K2 Thinking : 長時間推論と自律エージェントの最高峰 (2025年11月6日リリース)

2.1 アーキテクチャと革新的設計思想

Kimi K2 Thinkingは、中国のMoonshot AIが開発した**推論特化型大規模言語モデル**であり、「Model as Agent（モデル自体がエージェントとして機能）」という設計思想を体现している。アーキテクチャは**Mixture-of-Experts (MoE) を採用し、総パラメータ数は約1兆 (1T) という巨大な規模を誇るが、**推論時には約32億パラメータ (32B)** **のみが活性化される効率的な設計となっている。[\[3\]](#) [\[20\]](#) [\[21\]](#) [\[22\]](#)

最大の技術的特徴は、**「交互推論 (Interleaved Reasoning)」**と呼ばれる独自のメカニズムである。従来の推論モデルが「思考→最終回答」という一方向のプロセスをたどるのに対し、Kimi K2 Thinkingは「思考→ツール呼び出し→結果検証→さらなる思考」という循環的なプロセスを繰り返することで、より深く正確な推論を実現している。この設計により、**200~300回の連続的なツール呼び出し**を人間の介入なしに自律的に実行できる能力を持つ。[\[20\]](#) [\[23\]](#) [\[22\]](#) [\[24\]](#) [\[3\]](#)

コンテキストウィンドウは**256,000トークン (256K)**をサポートし、長大なドキュメント、複数ソースの研究資料、拡張された推論チェーンを一度に処理できる。さらに、**ネイティブINT4量子化 (Quantization-Aware Training)** を採用しており、**推論速度を従来比で約2倍向上させながら、性能の劣化をほぼゼロに抑える**ことに成功している。すべてのベンチマーク結果はINT4精度での測定値として報告されており、実運用環境での性能を正確に反映している。[\[21\]](#) [\[23\]](#) [\[22\]](#) [\[25\]](#) [\[3\]](#) [\[20\]](#)

2.2 ベンチマーク性能：多数の指標でSOTA達成

Kimi K2 Thinkingは複数の権威あるベンチマークで**State-of-the-Art (SOTA)**を達成し、特にクローズドソースモデルであるGPT-5やClaude Sonnet 4.5を上回る成績を記録している。[\[22\]](#) [\[26\]](#) [\[3\]](#) [\[20\]](#)

Humanity's Last Exam (HLE) では、ツール使用が許可される条件下で44.9%という正解率を達成し、これは同条件におけるSOTAスコアである。HLEは専門家レベルの知識と推論能力を問う極めて難易度の高いベンチマークであり、この成績は博士課程レベルの複雑な問題に対応できる能力を示している。[\[27\]](#) [\[3\]](#) [\[20\]](#)

数学的推論能力では驚異的な成績を記録しており、**AIME 2025** (米国数学招待試験) でPythonツール使用時に**99.1%**、ツールなしでも**94.5%の正解率を達成した**。さらに、ハーバード-MIT数学トーナメントを模した**HMMT 2025**でも、Pythonツール使用時に**95.1%**という極めて高いスコアを記録している。これらの成績は、競技数学レベルの問題解決能力を実証するものである。[\[23\]](#) [\[28\]](#) [\[20\]](#)

エージェント能力を測定する**BrowseComp**では**60.2%のスコアで首位に立ち、人間の平均的な成績29.2%を大きく上回っている。このベンチマークは、ウェブブラウジングを通じて難解な情報を自律的に収集・検証する能力を評価するものであり、実世界の研究タスクにおける有用性を示している。さらに、 **τ^2 -Bench Telecom**では93%**という驚異的なスコアを記録し、先行するK2 Instructの73%から大幅に向上している。[\[3\]](#) [\[20\]](#) [\[23\]](#) [\[22\]](#)

コーディング能力でも高い実力を発揮しており、**SWE-bench Verified**で71.3%、**LiveCodeBench v6**で**83.1%**を記録している。これらの成績は、実際のGitHub issueの解決や競技プログラミングにおける高い実用性を示すものである。[\[20\]](#) [\[23\]](#) [\[22\]](#)

2.3 推論プロセスの透明性と実用仕様

Kimi K2 Thinkingの独自性は、**推論過程の完全な透明性**にある。APIレスポンスには通常のcontentフィールド（最終回答）に加えて、reasoningフィールドが含まれており、モデルがどのような思考プロセスを経て結論に至ったかを詳細に確認できる。これにより、特に法務、医療、金融などの専門分野において、AIの判断根拠を検証・説明する必要がある場合に極めて有用である。[\[29\]](#) [\[23\]](#) [\[22\]](#) [\[20\]](#)

公式が公開した事例では、博士課程レベルの数学問題を**23回の推論とツール呼び出し**を通じて解決する過程が示されており、「仮説の提案→証拠の検証→推論の改良」という動的なサイクルを何百回も繰り返す能力が実証されている。このような長期的な計画と自律的な探索能力により、曖昧でオープンエンドな問題を明確で実行可能なサブタスクに分解できることが確認されている。[\[3\]](#) [\[20\]](#)

API価格は、標準版 (kimi-k2-thinking) が入力トークン**100万あたり0.60ドル**（キャッシュヒット時は0.15ドル）、出力トークン**100万あたり2.50ドル**となっている。高速版 (kimi-k2-thinking-turbo) は入力1.15ドル、出力8.00ドルと高価だが、より低レイテンシーでの推論が可能である。Claude 4.5 Sonnetの入力3ドル/出力15ドルと比較すると、Kimi K2 Thinkingは**約1/6のコスト**で同等以上の性能を提供している。[\[30\]](#) [\[31\]](#)

推論速度は典型的に**8~25秒**の範囲であり、単純なクエリでは高速だが複雑なタスクでは時間をかけて深く思考する。INT4量子化により、従来のFP8/BF16と比較して推論速度が約2倍向上しており、メモリフットプリントも大幅に削減されている（約600GBに対し、従来は1TB以上）。[\[25\]](#) [\[32\]](#) [\[22\]](#) [\[20\]](#)

2.4 適用領域と制約事項

Kimi K2 Thinkingが最も威力を発揮するのは、**長時間にわたる自律的な研究・分析タスク**である。学術研究、特許分析、市場調査、技術的な問題解決など、複数のソースから情報を収集し、仮説を検証し、論理的な結論を導く必要がある場合に最適である。実際に、複雑な戦略プランの立案やマーケティング施策の提案などで高い推論能力を発揮することが実証されている。[\[21\]](#) [\[23\]](#) [\[22\]](#) [\[20\]](#) [\[3\]](#)

一方で、高いハードウェア要件が制約となっており、INT4量子化を使用してもモデルの実行には相当なGPUメモリが必要である。ローカル環境での実行を検討する場合、複数の高性能GPUが必要となり、初期投資が大きくなる。また、創作的な文章生成においては「機械的に感じられる」という指摘があり、純粋な文学的創作には向かない可能性がある。[\[33\]](#) [\[22\]](#) [\[20\]](#)

商用利用においては、**Modified MITライセンス**の下で提供されており、大規模な商用利用には追加のライセンス要件が適用される場合がある。長いコンテキストと多数のツール呼び出しを伴うエージェント型ワークフローでは、トークン消費量が増加し、コストが予想以上に膨らむ可能性があるため、適切な予算管理が必要である。[\[31\]](#) [\[22\]](#) [\[30\]](#)

3. ERNIE-4.5-VL-28B-A3B-Thinking : マルチモーダル推論の新境地 (2025年11月11日リリース)

3.1 マルチモーダルアーキテクチャの革新

ERNIE-4.5-VL-28B-A3B-Thinkingは、中国Baiduが開発したマルチモーダル（視覚-言語）推論モデルであり、テキストと画像の両方を高度に理解・推論する能力を持つ。アーキテクチャは**Mixture-of-Experts (MoE) を採用し、総パラメータ数は28億（28B）だが、推論時に活性化されるのはわずか3億パラメータ（3B）**という極めて軽量な設計である。[\[4\]](#) [\[34\]](#) [\[35\]](#) [\[36\]](#)

この軽量設計にもかかわらず、複数のベンチマークで業界トップクラスのフラッグシップモデルに匹敵または凌駕する性能を達成している点が最大の革新である。公式発表によれば、従来の同等規模フルパラメータモデルと比較して推論速度が2~3倍高速であり、メモリフットプリントも大幅に削減されている。[\[34\]](#) [\[37\]](#) [\[36\]](#) [\[4\]](#)

技術的な特徴として、モデルは広範な中間学習段階（mid-training phase）を通じて、膨大な高品質視覚-言語推論データを吸収している。この大規模学習プロセスにより、視覚とテキストのモダリティ間の深い意味的アライメントが実現され、微妙な視覚-テキスト推論において前例のない能力を発揮する。[\[35\]](#) [\[38\]](#) [\[4\]](#) [\[34\]](#)

さらに、**GSPOとIcePop戦略を統合した最先端のマルチモーダル強化学習**を採用しており、動的難易度サンプリングと組み合わせることで、MoEトレーニングの安定性と学習効率を大幅に向上させている。[\[38\]](#) [\[4\]](#) [\[34\]](#) [\[35\]](#)

3.2 6つの主要能力と「画像と共に思考」機能

ERNIE-4.5-VL-28B-A3B-Thinkingは、6つの主要能力で卓越した性能を発揮する：[\[36\]](#) [\[4\]](#) [\[34\]](#) [\[35\]](#)

視覚推論 (Visual Reasoning)：大規模強化学習により強化され、複雑な視覚タスクにおける多段階推論、チャート分析、因果関係推論で例外的な能力を実証している。複雑な統計図表の分析や画像内の因果関係理解において正確な分析結果を提供する。[\[4\]](#) [\[34\]](#) [\[36\]](#)

STEM推論 (STEM Reasoning)：強力な視覚能力を活用し、写真からの数学問題解決、物理式の認識と計算、幾何図形分析などのSTEMタスクで飛躍的な性能向上を達成している。写真から直接数式や幾何図形を認識し、正確な計算と推論を実行できる。[\[34\]](#) [\[36\]](#) [\[4\]](#)

視覚的グラウンドィング (Visual Grounding)：コミュニティの強い要望に応え、より正確なオブジェクト位置特定と柔軟な指示実行能力を大幅に強化している。複雑な産業シナリオにおいて、言語プロンプトと視覚領域を正確にリンクし、構造化された出力を生成できる。[\[35\]](#) [\[36\]](#) [\[4\]](#) [\[34\]](#)

「画像と共に思考 (Thinking with Images)」：最も革新的な機能であり、画像ズームや画像検索などのツールと組み合わせることで、細粒度の詳細処理とロングテール視覚知識の取り扱い能力を劇的に向上させている。例えば、画像内の特定領域を自律的にズームインして詳細を確認し、看板の文字などを正確に識別できる。[\[4\]](#) [\[34\]](#) [\[35\]](#)

ツール呼び出し (Tool Calling)：画像検索、ズームなどのツールをネイティブにサポートし、視覚推論プロセスを動的に強化できる。これにより、内部知識だけでは不十分な場合に外部ツールを活用して情報を補完できる。[\[39\]](#) [\[36\]](#) [\[34\]](#) [\[4\]](#)

動画理解 (Video Understanding) : 強力な時系列認識能力を持ち、動画イベントの位置特定や動画内容の理解において優れた性能を発揮する。 [37] [40] [4]

3.3 性能指標と実用的デプロイメント

ERNIE-4.5-VL-28B-A3B-Thinkingは、複数のマルチモーダルベンチマークで優れた性能を示している。公式ベンチマークによれば、わずか3Bのアクティブパラメータで、**GPT-5-High**や**Gemini 2.5 Pro**を一部のテストで上回る性能を達成している。 [7] [41]

特筆すべきは、そのパラメータ効率の高さである。3Bアクティブパラメータのみでトップクラスのモデル性能を達成し、推論コストを**50%以上削減**している。推論速度は従来モデルの2~3倍に向上しており、メモリフットプリントも大幅に削減されている。単一推論レイテンシは200~500ms (入力長に依存)、スループットは20~50リクエスト/秒 (vLLM、単一A100使用時) という実用的な性能を提供する。 [37] [36] [4]

ライセンスは**Apache 2.0**であり、商用利用に制限がない完全なオープンソースモデルとして提供されている。Transformers、vLLM、FastDeployなど複数の推論フレームワークをサポートしており、企業環境での導入も容易である。 [36] [37] [4]

API価格は、入力トークンが**100万あたり0.42ドル**、出力トークンが**100万あたり1.25ドル**となっており、GPT-5と比較して入力は約1/3、出力は約1/8のコストである。この価格設定により、マルチモーダルAI機能を低コストで利用できる選択肢を提供している。 [42] [43]

3.4 適用領域と技術的制約

ERNIE-4.5-VL-28B-A3B-Thinkingが最も威力を発揮するのは、**視覚情報とテキスト情報を統合的に処理する必要があるタスク**である。具体的には、工業品質検査、医療画像診断、自動運転の環境認識、ロボットの視覚ナビゲーション、学術文献の図表分析、教育支援ツールなどが挙げられる。 [39] [7] [36] [4]

特に、中国語環境での最適化が進んでおり、中国語テキストと画像の組み合わせにおいて特に高い性能を発揮する。これは、Baiduが中国市場に深く根ざした企業であることを反映している。 [37] [4]

一方で、技術的な制約も存在する。モデルのロードには**80GB以上のGPUメモリ**が必要であり、単一のハイエンドGPU (A100やH100など) が必須となる。推論時のメモリ使用量は従来モデルより少ないものの、デプロイメントには相応のインフラ投資が必要である。 [36] [4] [37]

また、ベンチマーク情報が他の主要モデルと比較して限定的であり、特に英語圏の標準ベンチマークでの詳細な比較データが不足している。このため、グローバル市場での競争力を客観的に評価することがやや困難である。さらに、主に中国語に最適化されているため、英語や他の言語でのパフォーマンスがどの程度かは、さらなる検証が必要である。 [44] [7] [4] [37]

4. GPT-5.1：適応推論とユーザビリティの進化 (2025年11月12日リリース)

4.1 GPT-5からの主要な改良点

GPT-5.1は、OpenAIが2025年8月にリリースしたGPT-5の「アップグレード版」として位置づけられ、ユーザーエクスペリエンスと実用性の大幅な改善に焦点を当てている。最大の変更点は、モデルが**「Instant」と「Thinking」の2つのモードに分割され、タスクの複雑さに応じて自動的に適切なモードを選択する適応推論 (Adaptive Reasoning) **機能を搭載したことである。[\[5\]](#) [\[45\]](#) [\[46\]](#) [\[8\]](#) [\[9\]](#)

GPT-5.1 Instantは、速度と低レイテンシーを重視し、短いトランザクション型プロンプトや即座の出力が必要な自動化タスクに最適化されている。一方、**GPT-5.1 Thinking**は、複雑なワークフロー、データ解釈、コンテンツ企画など、深い推論リソースを必要とするタスクに特化している。この二層構造により、ユーザーは時間とコストの両方を最適化できる柔軟性を得ている。[\[45\]](#) [\[46\]](#) [\[9\]](#) [\[5\]](#)

GPT-5が「冷たい」「ロボット的」という評価を受けたことへの反省から、GPT-5.1は既定でより温かく、会話的なトーンを持つように調整されている。OpenAIは「楽しく会話できるAI」を目指し、有用性を保ちながら遊び心のあるモデルを実現したと述べている。さらに、ユーザーが「プロフェッショナル」「フレンドリー」「効率的」などのトーンプリセットを選択できるパーソナライゼーション機能も拡充されている。[\[8\]](#) [\[9\]](#) [\[47\]](#) [\[48\]](#) [\[5\]](#)

4.2 適応推論システムとベンチマーク性能

GPT-5.1の適応推論システムは、タスクの複雑さを動的に判断し、思考時間を調整する革新的な仕組みである。単純なクエリには最小限の計算リソースを割り当て、複雑な推論タスクには追加の分析層を提供する。OpenAIの内部ベンチマークによれば、GPT-5と比較して、最も簡単なタスクでは約2倍高速、**最も難しいタスクでは約2倍遅く（より徹底的）**動作する。[\[46\]](#) [\[9\]](#) [\[45\]](#)

代表的なChatGPTタスクの分布において、10パーセンタイル（単純タスク）では生成トークンが57%削減され、2倍高速な応答を実現している。一方、90パーセンタイル（複雑タスク）では71%多くのトークンを生成し、より徹底的な分析を行う。中央値（50パーセンタイル）では処理時間に変化がなく、全体として平均トークン消費量を23%削減しながら、複雑タスクの精度を18%向上させている。[\[45\]](#)

ベンチマーク性能では、**AIME 2025**（高度な数学問題）で**94.6%の精度を達成し、**GPT-4o**を8.4ポイント上回っている。**SWE-bench Verified**（実世界のコーディング課題）では74.9%のスコアを記録し、**GPT-4 Turbo**と比較して72.7%の改善を示している。**MMMU**（マルチモーダル理解）では84.2%**を達成し、最先端モデルの中でもトップ3にランクインしている。[\[49\]](#) [\[45\]](#)

コンテキストウィンドウは**400,000**トークンをサポートし、具体的には入力が**272,000**トークン、出力が**128,000**トークンという広大なキャパシティを持つ。これはGPT-4 Turboの128Kの4倍に相当し、コードベース全体や長文ドキュメントの比較分析を可能にしている。[\[50\]](#) [\[45\]](#)

4.3 API価格と実用的改善

GPT-5.1のAPI価格は、入力トークンが**100万あたり1.25**ドル、出力トークンが**100万あたり10.00**ドルとなっている。これは調査対象の4モデルの中で最も高価格帯であり、Claude Sonnet 4.5やKimi K2 Thinkingと比較して数倍のコストとなる。ただし、キャッシュ入力を活用すれば入力コストを**100万あたり0.125**ドルに削減でき、繰り返し使用される長いプロンプトコンテキストに対してコスト最適化が可能である。[\[51\]](#) [\[45\]](#)

実用的な改善点として、**指示追従能力の向上**が挙げられる。GPT-5.1は、「正確に50単語で書く」「表形式で出力する」といった制約をより確実に守るようになり、曖昧な表現に対する解釈精度も向上している。この改善は適応推論システムにより、プロンプトを処理する前により多くの時間をかけて微妙な制約を捉えることで実現されている。[\[9\]](#) [\[5\]](#) [\[45\]](#)

また、**専門用語の使用が減少し**、複雑な概念をより平易な言葉で説明する能力が向上している。これにより、技術的な専門知識を持たないユーザーでもAIの応答を理解しやすくなっている。幻覚（事実誤認）も大幅に削減されており、GPT-4oと比較して**45%少ない事実エラー**を生成することが報告されている。[\[48\]](#) [\[5\]](#) [\[8\]](#) [\[45\]](#)

4.4 エンタープライズ統合と実用性

GPT-5.1は、エンタープライズ環境での統合を強く意識した設計となっている。CRM更新、AIチャットアシスタント、ワークフロー自動化ツールなどのシステムにおいて、Instantモードを使用することで応答遅延を大幅に削減し、ユーザー エクスペリエンスとシステム効率の両方を改善できる。[\[46\]](#) [\[45\]](#)

一方、Thinkingモードはエンタープライズロジック、多段階推論、データ検証、コンプライアンスレビューなどの重要なプロセスに適している。これらの重いプロセスを軽量なプロセスから分離することで、精度と予算配分の両方をより細かく制御できるようになっている。[\[46\]](#)

OpenAIは、Enterprise、Business、Edu、Proプランのユーザーに対して、API呼び出し制限を**50%増加させ**、優先速度とリソース配分を提供している。さらに、使用率が90%に達すると、システムが自動的にGPT-5.1 Miniへの切り替えを推奨し、プロジェクトの進行が制限されないように配慮している。[\[52\]](#)

ただし、コストの高さは依然として大きな障壁であり、**トークン消費量の多いアプリケーションでは月額コストが数千ドルに達する可能性がある**。特に、Thinkingモードは複雑なタスクで多くのトークンを消費するため、過度に使用するとコストが急増するリスクがある。このため、OpenAIはタスクの複雑さに応じてInstantとThinkingを適切にルーティングし、コスト効率を最大化することを推奨している。[\[45\]](#) [\[46\]](#)

5. 4 モデルの総合比較：得意分野とポジショニング

5.1 コーディング・開発支援：MiniMax-M2の優位性

コーディングと開発支援のユースケースでは、**MiniMax-M2が最もバランスの取れた選択肢となる**。SWE-bench Verifiedで69.4%という高いスコアを記録し、実務的なプログラミング課題における強さを実証している。さらに重要なのは、その圧倒的なコスト効率と推論速度であり、Claude Sonnet 4.5の約8%のコストで2倍の推論速度を実現している。[\[53\]](#) [\[11\]](#) [\[33\]](#) [\[1\]](#) [\[16\]](#)

開発者コミュニティの報告では、「書く→編集→再生成→改良」のサイクルにおいて、MiniMax-M2の高速性が生産性を劇的に向上させるという評価が多い。特に、スタートアップや中小企業のように、予算制約がある中で高品質なAI支援を求める場合、MiniMax-M2は理想的な選択肢となる。[\[33\]](#) [\[2\]](#) [\[16\]](#)

ただし、複雑なリポジトリレベルのバグ修正や、極めて高度なアーキテクチャ設計においては、GPT-5.1やClaude Sonnet 4.5がわずかに優位に立つ場合がある。このため、**70%の標準的なコーデ**

イングタスクにMiniMax-M2を使用し、30%の複雑なタスクにGPT-5.1を使用するハイブリッド戦略が最もコスト効率が高いとされる。[32] [33] [45]

5.2 数学・競技プログラミング：Kimi K2 Thinkingの圧倒的強さ

数学的推論と競技プログラミングにおいては、**Kimi K2 Thinking**が他を圧倒する性能を示している。AIME 2025で99.1%（Python使用時）、HMMT 2025で95.1%（Python使用時）という驚異的なスコアは、競合モデルを大きく引き離している。LiveCodeBench v6でも83.1%という高スコアを記録し、競技プログラミングレベルのアルゴリズム問題に対する強さを実証している。[28] [23] [22] [20]

これらの成績は、Kimi K2 Thinkingの「交互推論」メカニズムが、数学的証明やアルゴリズム設計のような段階的思考をするタスクに極めて適していることを示している。200～300回の連続的なツール呼び出しを通じて、複雑な問題を小さなサブ問題に分解し、各ステップで検証を行いながら解決する能力は、人間の数学者のアプローチに近い。[24] [23] [20] [3]

ただし、この高性能には相応のコストが伴う。推論時間が8～25秒と長く、リアルタイムのコーディング支援には向かない。また、複雑な数学問題ではトークン消費量が大幅に増加し、API使用コストが予想以上に高くなる可能性がある。このため、数学研究、アルゴリズム開発、高度な技術問題解決といった、正確性と深い推論が最優先されるユースケースに限定して使用することが推奨される。[30] [32] [31] [20] [33]

5.3 マルチモーダル理解：ERNIE-4.5-VL-28Bの独自性

視覚情報とテキスト情報を統合的に処理する必要がある場合、**ERNIE-4.5-VL-28B-A3B-Thinking**は唯一の真のマルチモーダル推論モデルとして独自のポジションを占める。画像からの数学問題解決、工業品質検査、医療画像診断、自動運転の環境認識など、視覚情報が重要な役割を果たすタスクにおいて、他のテキスト専用モデルでは不可能な能力を提供する。[7] [34] [4] [36]

特に、「画像と共に思考」機能は革新的であり、画像内の特定領域を自律的にズームインして詳細を確認したり、外部の画像検索ツールを呼び出して追加情報を収集したりすることができる。これにより、単一の画像から得られる情報を最大限に活用し、高精度な分析を実現している。[34] [35] [4] [36]

コスト面でも魅力的であり、入力0.42ドル/出力1.25ドル（100万トークンあたり）という価格設定は、GPT-5.1の約1/3～1/8のコストである。さらに、わずか3Bのアクティブパラメータで推論速度が従来モデルの2～3倍という効率性は、大規模な画像処理タスクにおいて大きな優位性となる。[43] [42] [4] [37] [36]

ただし、主に中国語環境での最適化が進んでいるため、英語や他の言語での性能については追加検証が必要である。また、80GB以上のGPUメモリが必要というハードウェア要件は、ローカルデプロイメントのハードルとなる。[4] [37] [36]

5.4 会話・カスタマーサポート：GPT-5.1の温かみと柔軟性

会話型AIやカスタマーサポートのユースケースでは、**GPT-5.1の温かいトーンと適応推論が大きな強みとなる**。GPT-5が「冷たい」という批判を受けたことへの反省から、GPT-5.1は既定でより親しみやすく、会話を楽しめるモデルとして再設計されている。[5] [8] [9] [48]

適応推論により、単純な問い合わせには迅速に応答し（Instantモード）、複雑な問題にはじっくり考えて対応する（Thinkingモード）という柔軟な対応が可能である。カスタマーサポート自動化の実験

では、単純な問い合わせの76%を4~7秒で正確に解決し、顧客満足度91%を達成している。[\[32\]](#) [\[5\]](#) [\[45\]](#) [\[46\]](#)

さらに、「プロフェッショナル」「フレンドリー」「効率的」などのトーンプリセットにより、ブランドガイドラインに沿った一貫したコミュニケーションを実現できる。これは、企業の顧客対応において極めて重要な要素である。[\[47\]](#) [\[9\]](#) [\[5\]](#) [\[46\]](#)

ただし、高コストが最大の障壁となる。出力トークンが100万あたり10.00ドルという価格設定は、大量の顧客対応を行う場合に月額コストが数万ドルに達する可能性がある。このため、**高価値顧客向けの対応にGPT-5.1を使用し、標準的な問い合わせにはMiniMax-M2を使用する階層的戦略**が推奨される。[\[51\]](#) [\[33\]](#) [\[32\]](#) [\[45\]](#)

5.5 エンタープライズ統合と長期運用

エンタープライズ環境での統合と長期運用を考慮する場合、**GPT-5.1が最も成熟したエコシステムと信頼性**を提供する。OpenAIは長年にわたってエンタープライズ顧客との関係を構築しており、SLA、優先処理、専用サポートなどの包括的なサービスを提供している。[\[51\]](#) [\[45\]](#) [\[46\]](#)

CRM、ワークフロー自動化、ドキュメント管理システムなど、既存のエンタープライズアプリケーションとの統合においても、GPT-5.1は豊富なツールとライブラリを持つ。さらに、OpenAIの継続的な投資により、モデルの長期的なサポートと更新が保証されている。[\[9\]](#) [\[45\]](#) [\[46\]](#)

一方、**中国企業が開発した3モデル（MiniMax-M2、Kimi K2 Thinking、ERNIE-4.5-VL-28B）**は、コスト効率において圧倒的な優位性を持つ。特に、大規模な自動化タスクやバッチ処理を行う場合、これらのモデルを使用することで年間数十万ドルのコスト削減が可能となる。[\[54\]](#) [\[11\]](#) [\[42\]](#) [\[1\]](#) [\[16\]](#)

現実的な戦略としては、**コアビジネスロジックや顧客対応にGPT-5.1を使用し、バックグラウンドの分析やデータ処理に中国製モデルを使用するハイブリッドアプローチ**が、コストと性能のバランスを最適化する。このアプローチにより、エンタープライズの信頼性を維持しながら、運用コストを50~70%削減できる可能性がある。[\[16\]](#) [\[33\]](#) [\[32\]](#)

6. 2025年11月時点の生成AI技術動向と今後の展望

6.1 オープンソースとクローズドソースのギャップ縮小

2025年11月の4モデルリリースは、**オープンソースとクローズドソースモデルのギャップが急速に縮小**していることを明確に示している。Artificial Analysisのデータによれば、最高のオープンソースモデル（MiniMax-M2、品質スコア61）と最高のプロプライエタリモデル（GPT-5、推定68）のギャップは現在約7ポイントであり、昨年の約18ポイントから大幅に縮小している。線形外挿によれば、このギャップは2026年中頃には実質的に消滅する可能性がある。[\[55\]](#) [\[6\]](#) [\[22\]](#) [\[2\]](#) [\[3\]](#)

特に、**Kimi K2 ThinkingがHLE、BrowseComp、数学ベンチマークでGPT-5やClaude Sonnet 4.5を上回った**事実は、中国勢の技術的キャッチアップが完了したことを示す象徴的な出来事である。これは、AI開発の中心が徐々に米国から多極化している構造的变化を反映している。[\[26\]](#) [\[6\]](#) [\[22\]](#) [\[3\]](#)

6.2 効率性の時代：MoEアーキテクチャの主流化

調査対象の4モデルのうち3モデル（MiniMax-M2、Kimi K2 Thinking、ERNIE-4.5-VL-28B）が **Mixture-of-Experts (MoE)** アーキテクチャを採用していることは、効率性が最優先課題となっていることを示している。MoEにより、「大規模な知識容量を維持しながら、推論時のコストと速度を最適化する」という理想的なバランスが実現されている。[\[22\]](#) [\[1\]](#) [\[2\]](#) [\[33\]](#) [\[3\]](#) [\[4\]](#)

特に注目すべきは、アクティブパラメータの劇的な縮小である。ERNIE-4.5-VL-28Bはわずか3Bのアクティブパラメータでトップクラスの性能を達成し、Kimi K2 Thinkingは1兆パラメータのモデルを32Bのアクティブで運用している。この「スパース化」のトレンドは、今後さらに加速すると予想される。[\[12\]](#) [\[22\]](#) [\[3\]](#) [\[4\]](#)

6.3 推論特化型モデルの台頭：Test-Time Scalingの重要性

Kimi K2 ThinkingとGPT-5.1の両方が**「推論時間を延長することで精度を向上させる」というTest-Time Scalingアプローチ**を採用している点は、今後の技術発展の方向性を示している。従来の「モデルを大きくすれば性能が上がる」というパラダイムから、「推論時に時間をかけければ精度が上がる」という新しいパラダイムへの移行が進んでいる。[\[20\]](#) [\[22\]](#) [\[9\]](#) [\[3\]](#) [\[45\]](#)

この変化により、開発者は**「速度と精度のトレードオフを動的に調整する」**能力を得ている。単純なタスクには高速モードを使用し、複雑なタスクには深い思考モードを使用するという柔軟な戦略が、コストパフォーマンスを最大化する新しいベストプラクティスとなりつつある。[\[9\]](#) [\[33\]](#) [\[32\]](#) [\[45\]](#) [\[46\]](#)

6.4 マルチモーダルの深化：視覚推論の次のフロンティア

ERNIE-4.5-VL-28Bが示すように、テキストと視覚の統合推論能力が次の重要なフロンティアとなっている。単に「画像を認識する」だけでなく、「画像と共に思考する」「画像情報に基づいて推論する」能力が、実世界のタスクにおいて決定的に重要な場面が増えている。[\[35\]](#) [\[7\]](#) [\[34\]](#) [\[36\]](#) [\[4\]](#)

工業、医療、自動運転、教育などの分野では、視覚情報が意思決定の中核を占めるため、真のマルチモーダル推論能力は今後さらに需要が高まると予想される。特に、「画像と共に思考」や「視覚的グラウンディング」のような、視覚情報を動的に操作しながら推論するアプローチは、次世代AIシステムの標準機能となる可能性が高い。[\[7\]](#) [\[34\]](#) [\[35\]](#) [\[36\]](#) [\[4\]](#)

6.5 コスト競争の激化：価格破壊がもたらす民主化

MiniMax-M2の「Claude比8%のコスト」やKimi K2 Thinkingの「Claude比約1/6のコスト」という価格設定は、AI市場における価格破壊が本格化していることを示している。この価格競争は、中国企業による積極的な市場参入と、DeepSeekが2024年5月に引き起こした「第一次価格戦争」の延長線上にある。[\[11\]](#) [\[54\]](#) [\[31\]](#) [\[1\]](#) [\[30\]](#)

この価格破壊は、AIの民主化を大きく促進している。従来は大企業しか利用できなかった高性能AIが、スタートアップや中小企業、さらには個人開発者にも手の届く価格となり、イノベーションの裾野が広がっている。特に、コスト制約のある新興市場や発展途上国において、これらの低価格モデルがAI活用の障壁を大幅に下げることが期待される。[\[54\]](#) [\[1\]](#) [\[2\]](#) [\[16\]](#)

一方で、OpenAIのような高価格帯のプロバイダーは、エンタープライズ向けの信頼性、サポート、統合容易性を差別化要因として強調している。この「低価格・高効率 vs 高価格・高信頼性」という

市場の二極化は、今後さらに明確になると予想される。[33] [45] [46] [51]

まとめ

2025年10月下旬から11月中旬にかけてリリースされた4つの最先端生成AIモデル——MiniMax-M2、Kimi K2 Thinking、ERNIE-4.5-VL-28B-A3B-Thinking、GPT-5.1——は、それぞれ異なる技術的アプローチと市場ポジショニングを持ち、生成AI業界の多様化と成熟化を象徴している。[1] [5] [3] [4]

MiniMax-M2は、コスト効率と推論速度のバランスにおいて革新的であり、開発者やスタートアップにとって理想的な選択肢となる。**Kimi K2 Thinking**は、長時間推論と自律エージェント能力において業界最高峰の性能を示し、数学や複雑な問題解決において無類の強さを發揮する。**ERNIE-4.5-VL-28B-A3B-Thinking**は、真のマルチモーダル推論能力を持つ唯一のモデルとして独自のポジションを占め、視覚情報が重要な役割を果たすタスクにおいて不可欠である。**GPT-5.1**は、会話性とユーザビリティにおいて最も洗練されており、エンタープライズ統合と顧客対応において依然として優位性を持つ。[23] [2] [11] [22] [5] [1] [16] [3] [20] [34] [36] [7] [45] [46] [9] [4]

これら4モデルの比較分析から明らかになったのは、「単一の最良モデル」という概念がもはや存在せず、ユースケースに応じて最適なモデルを選択する戦略的思考が不可欠となっている点である。さらに、複数モデルを組み合わせたハイブリッドアプローチが、コストと性能の最適なバランスを実現する現実的な解決策となっている。[53] [32] [33]

2025年11月時点において、生成AI技術は「オープンソースの台頭」「効率性の追求」「推論特化型の進化」「マルチモーダルの深化」「価格競争の激化」という5つの明確なトレンドを示している。これらのトレンドは、今後数年間にわたってAI業界の発展を方向づける重要な要因となるだろう。企業や開発者は、これらの急速な変化に適応し、自らのニーズに最も適したモデルを戦略的に選択・統合することで、AI技術の恩恵を最大化できる。[55] [2] [22] [32] [54] [3] [33] [4]

※

1. <https://note.com/yaandyu0423/n/n25fd31691043>
2. <https://venturebeat.com/ai/minimax-m2-is-the-new-king-of-open-source-langs-especially-for-agentic-tools/>
3. https://note.com/trans_n_ai/n/nfc9888d1be28
4. <https://zenn.dev/czmilo/articles/fcaa5dca44d10c>
5. <https://www.itmedia.co.jp/aiplus/articles/2511/13/news053.html>
6. <https://japan.zdnet.com/article/35240260/>
7. <https://www.artificialintelligence-news.com/news/baidu-ernie-multimodal-ai-gpt-and-gemini-benchmarks/>
8. <https://forest.watch.impress.co.jp/docs/news/2062843.html>
9. <https://www.datacamp.com/blog/gpt-5-1>
10. <https://techcommunity.microsoft.com/blog/azure-ai-foundry-blog/minimax-m2-the-open-source-innovator-in-coding-and-agentic-workflows-now-in-azur/4466045>
11. <https://technode.com/2025/10/28/minimax-releases-m2-open-source-model-offering-double-speed-at-8-of-claude-sonnets-price/>
12. <https://artificialanalysis.ai/articles/minimax-m2-benchmarks-and-analysis>
13. <https://openrouter.ai/minimax/minimax-m2>

14. <https://www.c-sharpcorner.com/article/minimax-m2-on-ollama-cloud-benchmark-leader-for-coding-and-agentic-workflows/>
15. <https://note.com/eiji71/n/n6d495c78a686>
16. <https://magichour.ai/blog/minimax-m2-vs-gpt-4o-vs-claude-35>
17. <https://binaryverseai.com/minimax-m2-review-setup-pricing-benchmarks-agent/>
18. <https://note.com/alvis8039/n/nf034b76188b5>
19. <https://zenn.dev/saan/articles/2db2c3b0939991>
20. <https://www.datacamp.com/tutorial/kimi-k2-thinking-guide>
21. <https://gai.workstyle-evolution.co.jp/2025/11/11/kimi-k2-thinking-1-trillion-parameter-open-source-ai-beats-gpt4-claude-sonnet-comprehensive-analysis/>
22. <https://felloai.com/it/2025/11/new-chinese-model-kimi-k2-thinking-ranks-1-in-multiple-benchmarks/>
23. <https://www.together.ai/models/kimi-k2-thinking>
24. <https://www.youtube.com/watch?v=ceglxJwDPyk>
25. <https://www.interconnects.ai/p/kimi-k2-thinking-what-it-means>
26. <https://36kr.jp/444000/>
27. <https://moonshotai.github.io/Kimi-K2/thinking.html>
28. <https://jobirun.com/moonshot-kimi-k2-thinking-deep-reasoning-model/>
29. https://oneword.co.jp/bignite/ai_news/china-kimi-k2-thinking-model-reasoning-ai-evolution/
30. <https://www.cometapi.com/how-to-use-kimi-k2-thinking-api-a-practical-guide/>
31. <https://apidog.com/blog/kimi-k2-thinking-api/>
32. <https://www.cursor-ide.com/blog/kimi-2-thinking-vs-gpt-5>
33. <https://go.lightnode.com/tech/minimax-m2-vs-glm4.6-vs-kimi-k2-thinking>
34. <https://huggingface.co/baidu/ERNIE-4.5-VL-28B-A3B-Thinking>
35. <https://ernie.baidu.com/blog/posts/ernie-4.5-vl-28b-a3b-thinking/>
36. <https://dev.to/czmilo/2025-complete-guide-in-depth-analysis-of-ernie-45-vl-28b-a3b-thinking-multimodal-ai-model-1mib>
37. <https://zenn.dev/qingwu/articles/0ee5df24da737e>
38. <https://aistudio.baidu.com/modelsdetail/39280/intro>
39. <https://www.marktechpost.com/2025/11/11/baidu-releases-ernie-4-5-vl-28b-a3b-thinking-an-open-source-and-compact-multimodal-reasoning-model-under-the-ernie-4-5-family/>
40. <https://www.youtube.com/watch?v=nbXgKPNpJJg>
41. <https://venturebeat.com/ai/baidu-unveils-proprietary-ernie-5-beating-gpt-5-performance-on-charts>
42. <https://www.labellerr.com/blog/baidu-launches-ernie-4-5-and-x1/>
43. <https://blog.galaxy.ai/compare/ernie-4-5-vl-424b-a47b-vs-gpt-5>
44. <https://news.aibase.com/ja/news/22760>
45. <https://www.cursor-ide.com/blog/gpt-51-vs-claude-45>
46. <https://scalevise.com/resources/gpt-5-1-new-features/>
47. <https://www.theverge.com/news/802653/openai-gpt-5-1-upgrade-personality-presets>
48. https://www.gizmodo.jp/2025/11/openai_chatgpt_gpt_5_1_released.html

49. https://note.com/r1250_gs/n/nee9c7c427f7f
50. https://www.vals.ai/models/openai_gpt-5.1-2025-11-13
51. <https://openai.com/api/pricing/>
52. <https://www.aibase.com/news/22651>
53. <https://vpssos.com/minimax-m2-vs-glm-4.6-vs-kimi-k2-thinking/>
54. <https://www.caixinglobal.com/2025-10-28/minimax-unveils-m2-model-to-compete-on-speed-and-cost-102376624.html>
55. https://www.reddit.com/r/LocalLLaMA/comments/1oihbtz/minimaxm2_cracks_top_10_overall_llms_production/
56. <https://www.rival.tips/compare/minimax-m2-free/kimi-k2-thinking>
57. <https://www.vellum.ai/llm-leaderboard>
58. <https://blog.galaxy.ai/compare/kimi-k2-vs-minimax-m2>
59. <https://venturebeat.com/ai/openai-reboots-chatgpt-experience-with-gpt-5-1-after-mixed-reviews-of-gpt-5>
60. <https://www.siliconflow.com/articles/en/best-LLMs-for-reasoning-tasks>
61. https://www.reddit.com/r/LocalLLaMA/comments/1ot3ueh/anyone_got_the_chance_to_compare_local_minimaxm2/
62. https://yiyan.baidu.com/blog/publication/ERNIE_Technical_Report.pdf
63. <https://openrouter.ai/models>
64. https://gigazine.net/gsc_news/en/20251028-minimax-m2-open-sourcing/
65. <https://cloud.google.com/vertex-ai/generative-ai/docs/maas/kimi/kimi-k2-thinking>
66. <https://weel.co.jp/media/tech/gpt-5-codex/>
67. <https://buildingclub.info/minimax-minimax-m2-cost-calculator-for-api-token-pricing/>
68. <https://openrouter.ai/moonshotai/kimi-k2-thinking>
69. https://www.reddit.com/r/LocalLLaMA/comments/1orrddh/minimax_m2_coding_plan_pricing_revealed/
70. <https://huggingface.co/moonshotai/Kimi-K2-Thinking>
71. <https://www.youtube.com/watch?v=07TnqXM-WUg>
72. <https://huggingface.co/blog/MiniMax-AI/why-did-m2-end-up-as-a-full-attention-model>
73. <https://www.youtube.com/watch?v=021Q7Q8XTPg>
74. <https://www.linkedin.com/pulse/fireworks-ai-october-2025-roundup-fireworks-ai-mnkqf>
75. https://www.linkedin.com/posts/sebastianraschka_i-just-saw-the-benchmarks-of-the-new-open-weight-activity-7388978810781827072-7Jgw
76. <https://news.aibase.com/news/22592>
77. <https://zenn.dev/beagle/scraps/96a556a831eb1a>
78. <https://github.com/MiniMax-AI/MiniMax-M2>
79. <https://artificialanalysis.ai/models/kimi-k2-thinking>
80. <https://www.cnbc.com/2025/11/06/alibaba-backed-moonshot-releases-new-ai-model-kimi-k2-thinking.html>
81. https://ledge.ai/articles/kimi_k2_thinking_open_source_release_2025

82. <https://venturebeat.com/ai/moonshots-kimi-k2-thinking-emerges-as-leading-open-source-ai-outperforming/>
83. <https://weel.co.jp/media/tech/kimi-k2-thinking/>
84. <https://innovatopia.jp/ai/ai-news/71214/>
85. https://www.reddit.com/r/LocalLLaMA/comments/1oqi4qp/my_handson_review_of_kimi_k2_thinking_the/
86. <https://staffing.archetyp.jp/magazine/kimi-k2-thinking/>
87. <https://innovatopia.jp/ai/ai-news/71466/>
88. <https://weel.co.jp/media/tech/ernie-4-5-vl-28b-a3b-thinking/>
89. <https://skywork.ai/blog/ja/models/baidu-ernie-4-5-vl-28b-a3b-free-chat-online/>
90. <https://innovatopia.jp/ai/ai-news/71772/>
91. <https://weel.co.jp/media/tech/ernie-5-0/>
92. <https://dera.ai/ja/news/9877bdd7-9f8f-36d1-53ef-23727c81adad>
93. <https://www.aibase.com/news/www.aibase.com/ja/news/22703>
94. <https://www.vietnam.vn/ja/ai-trung-quoc-danh-bai-gpt-5>
95. <https://xenospectrum.com/baidu-ernie-5-omnimodal-ai-beats-gpt5/>
96. <https://skywork.ai/skypage/en/chatgpt-5-1-features-benchmarks-future/1988849353132838912>
97. <https://note.com/npaka/n/nc00eea85734a>
98. <https://openai.com/index/gpt-5-system-card-addendum-gpt-5-1/>
99. <https://gihyo.jp/article/2025/11/gpt-5.1>
100. https://cdn.openai.com/pdf/4173ec8d-1229-47db-96de-06d87147e07e/5_1_system_card.pdf
101. <https://arstechnica.com/ai/2025/11/openai-walks-a-tricky-tightrope-with-gpt-5-1s-eight-new-personalities/>
102. <https://momo-gpt.com/column/chatgpt-5-1/>
103. <https://openai.com/index/gpt-5-1-for-developers/>
104. https://ledge.ai/articles/openai_gpt_5_1_release
105. https://www.perplexity.ai/page/openai-releases-gpt-5-1-api-wi-kkcfAIRIQ3KcXAPG2IQ_uw
106. <https://openai.com/ja-JP/index/introducing-gpt-5/>
107. <https://emma-benchmark.github.io>
108. <https://aclanthology.org/2025.findings-acl.1112.pdf>
109. <https://blog.galaxy.ai/compare/kimi-k2-thinking-vs-minimax-m2>
110. <https://blog.galaxy.ai/compare/ernie-4-5-vl-424b-a47b-vs-gpt-5-codex>
111. <https://arxiv.org/html/2509.14142v1>
112. <https://sourceforge.net/software/compare/Kimi-K2-Thinking-vs-MiniMax-M2/>
113. <https://skywork.ai/blog/ai-agent/gpt-5-1-thinking-ultimate-guide-complex-problem-solving/>
114. https://openaccess.thecvf.com/content/ICCV2025W/MARS2/papers/Xu_MARS2_2025_Challenge_on_Multimodal_Reasoning_Datasets_Methods_Results_Discussion_ICCVW_2025_paper.pdf