

GPT-5 vs OpenAI o3性能比較分析

Claude Opus 4.1

GPT-5は2025年8月にリリースされた統合型AIモデルで、OpenAI o3シリーズを大幅に上回る実用性能を示している。OpenAI数学、物理学、プログラミング分野で特に顕著な優位性を発揮し、50-80%少ないトークンでより高い精度を達成している。一方、o3は理論的推論において卓越した能力を持つものの、実用性とコスト効率の面で劣勢にある。

TechCrunch +3

モデルリリース時期と基本仕様

GPT-5の展開状況 GPT-5は2025年8月7日に正式リリースされ、ChatGPT、OpenAI API、GitHub Models Playgroundで利用可能となった。Botpress +4 統合アーキテクチャを採用し、**1.7兆パラメータ**の完全版に加え、Mini、Nano、Chatの各バリエントを提供している。OpenAI 272,000入力トークン、128,000出力トークンの処理能力を持つ。Simon Willison +3

OpenAI o3シリーズの段階的リリース o3シリーズは段階的にリリースされ、o3-miniが2025年1月31日、完全版o3が4月16日、o3-proが6月10日に展開された。Wikipedia +3 推論特化型アーキテクチャを採用し、200,000トークンのコンテキスト窓を持つ。「Artificial Analysis +3 私的思考連鎖」による深い推論機能が特徴である。Wikipedia +2

数学・物理学分野での性能比較

数学コンテスト・競技での成績

AIME（アメリカ数学招待試験）2025年版

- **GPT-5:** 94.6% (ツール無し)、OpenAI 99.6% (思考モード有効時) The Decoder +3
- **o3:** 88.9% (2025年版)、TechTarget OpenAI 96.7% (2024年版) OpenAI +4

この結果は、GPT-5が高度な数学問題において**6ポイント以上の優位性**を持つことを示している。

FrontierMath（研究レベル数学）

- **o3:** 25.2% (当初発表) DataCamp +4 → **10%** (独立検証結果)
- **GPT-5:** 具体的数値未公表だが、o3を上回る性能

物理・科学推論ベンチマーク GPQA Diamond (博士レベル物理・化学・生物学) arXiv
Papers with Code では：

- **GPT-5 Pro: 88.4%** (ツール無し)、89.4% (Pythonツール有り) Vellum OpenAI
- **o3:** 83.3% (一般版)、DataCamp 87.7% (一部評価) Wikipedia +3
- 人間の専門家 (博士号保持者) : 65-74% arXiv OpenReview

両モデルとも人間の専門家を大幅に上回る性能を示すが、GPT-5が一貫して優位性を保っている。

プログラミング・コーディング性能

ソフトウェア開発ベンチマーク

SWE-bench Verified (実世界ソフトウェア工学)

- **GPT-5: 74.9%** (新記録) OpenAI
- **o3:** 69.1%、o3-mini: 71.7% OpenAI +2
- 従来最高 (o1) : 48.9% Klu

Aider Polyglot (多言語コード編集)

- GPT-5: 88% (o3比でエラー率33%削減) (OpenAI)
- o3: GPT-5より低い成績

競技プログラミング

Codeforces評価

- o3: ELOレーティング**2,727** (99.7パーセンタイル、世界トップ200プログラマー相当)
(Medium +2)
- この分野では、o3がGPT-5を上回る例外的な性能を示している

化学・論理推論分野

化学推論能力

GPQA化学コンポーネント (arXiv) (Papers with Code)において、両モデルとも博士レベルの化学問題に対応可能だが、**GPT-5が総合的に優位性**を示している。化学オリンピック関連の具体的ベンチマークは限定的だが、両モデルとも博士レベルの理解を実証している。

抽象推論・論理

ARC-AGI (抽象推論チャレンジ)

- o3: 75.7% (低計算量)、**87.5%** (高計算量設定) (DataCamp +2)
- 人間の平均: 85%閾値 (DataCamp +2)
- GPT-5: 具体的数値未公表だが、全体的性能傾向から o3を上回ると予想

この分野では、o3が突出した性能を発揮し、AGIレベルの抽象推論能力に近づいている。

(TechCrunch +3)

技術アーキテクチャと改善点

GPT-5の技術革新

統合適応アーキテクチャ GPT-5は「リアルタイムルーター」システムを採用し、クエリの複雑さに応じて自動的に高速モード (gpt-5-main) と深い推論モード (gpt-5-thinking) を切り替える。 (OpenAI +2) この統合アプローチにより、専用モデル間の切り替えが不要となった。

(Simon Willison)

安全性向上 「セーフコンプリーション」訓練手法を導入し、二進拒否訓練に代わって安全制約内での有用性を最大化する。 (OpenAI +3) 結果として、**GPT-4oと比べて事実誤認が45%減少** (OpenAI) している。 (OpenAI)

o3の技術的特徴

推論特化アーキテクチャ o3は「熟考アライメント」手法を使用し、推論プロセス中に安全ポリシーについて明示的に思考する。 (OpenAI +2) テスト時探索アーキテクチャにより、最適解選択前に複数の推論経路を探索する。 (Adaline +2)

視覚推論統合 o3は推論モデルとして初めて、思考連鎖に画像を直接統合する機能を実装している。 (DataCamp)

実用性能とコスト効率

処理効率とハルシネーション

出力効率性

- **GPT-5**: o3比で**50-80%少ない出力トークン**で同等以上の性能 (OpenAI)
- **GPT-5**: 思考モード使用時の事実誤認率4.8% (TechCrunch) (OpenAI)
- **o3**: ハルシネーション率22% (TechCrunch) (OpenAI)

価格設定

API価格（100万トークンあたり）

- **GPT-5**: 入力\$1.25、出力\$10 (OpenAI) (Vellum)
- **o3**: 入力\$2、出力\$8 (80%価格削減後) (OpenAI Developer Community...)
- **o3-pro**: 入力\$20、出力\$80 (OpenAI Developer Community...)

実際の使用においてGPT-5は効率性により**50-80%のコスト削減**を実現している。

各モデルの強みと弱み

GPT-5の強み

1. **コーディング優秀性**: 開発者から「明確に最高のコーディングモデル」との評価 (OpenAI)
(Substack)
2. **ツール統合**: 並列ツール呼び出しと自律エージェント機能に優秀 (OpenAI)
3. **実用性**: 複雑なタスクを一回の試行で完成させる能力 (Substack)
4. **生産準備性**: プロトタイプではなく実用可能なコードを生成
5. **信頼性**: 思考モード使用時にo3と比べて**事実誤認が80%少ない** (OpenAI)

GPT-5の弱み

1. **創作執筆**: GPT-4.5と比べて創造的執筆能力が大幅に劣る (Substack)
2. **コンテンツ作成**: より「LinkedIn風の」機械的な応答傾向 (Substack)
3. **主観的タスク**: 創造性を要求される分野での劣勢

o3の強み

1. **学術推論**: 理論的・研究指向タスクに優秀
2. **深い思考**: 複雑な理論問題の徹底分析
3. **競技プログラミング**: Codeforces ELO 2,727の卓越した成績 (Medium +5)

o3の弱み

1. **実装ギャップ**: 作業ソリューションではなく計画提供の傾向 (Substack)
2. **速度**: 推論に高いレイテンシー
3. **コスト効率**: 同等性能でより高額
4. **ツール回復**: ツール呼び出し失敗からの回復能力が低い (Substack)

結論と推奨事項

GPT-5は実用アプリケーションにおいてo3を大幅に上回る性能を示している。特に**数学 (AIME 2025: 94.6% vs 88.9%)、コーディング (SWE-bench: 74.9% vs 69.1%)、医療応用 (HealthBench: 46.2% vs 31.6%) **で優位性が顕著である。 (OpenAI)

GPT-5使用推奨場面：生産ソフトウェア開発、コスト効率重視のAI統合、大規模コードベース作業、実用的技術問題解決 (OpenAI)

o3使用推奨場面：理論研究、予算制約の少ない専門STEM問題、最高品質の推論が必要な場面

実用AGI応用への意味ある前進として、GPT-5は統合体験と自動ルーティングにより、o3の手動モデル選択に存在した多くのユーザー摩擦点を解消している。早期アクセス開発者と企業ユーザーの合意として、GPT-5は特にソフトウェア工学において実用的AGI応用への意味ある一歩を示している（Substack）一方、o3はコストより絶対的推論品質が重要な研究・学術応用により適している。