

# 2026年 大学入学共通テストにおける生成AI の性能評価と教育・産業への多層的影響分析 ～GPT-5.2 Thinking、Gemini 3 Pro、Claude 4.5 Opusの 比較と「知能」の再定義～

Gemini 3 pro

## 1. 序論：標準化試験における特異点の出現とパラダイムシフト

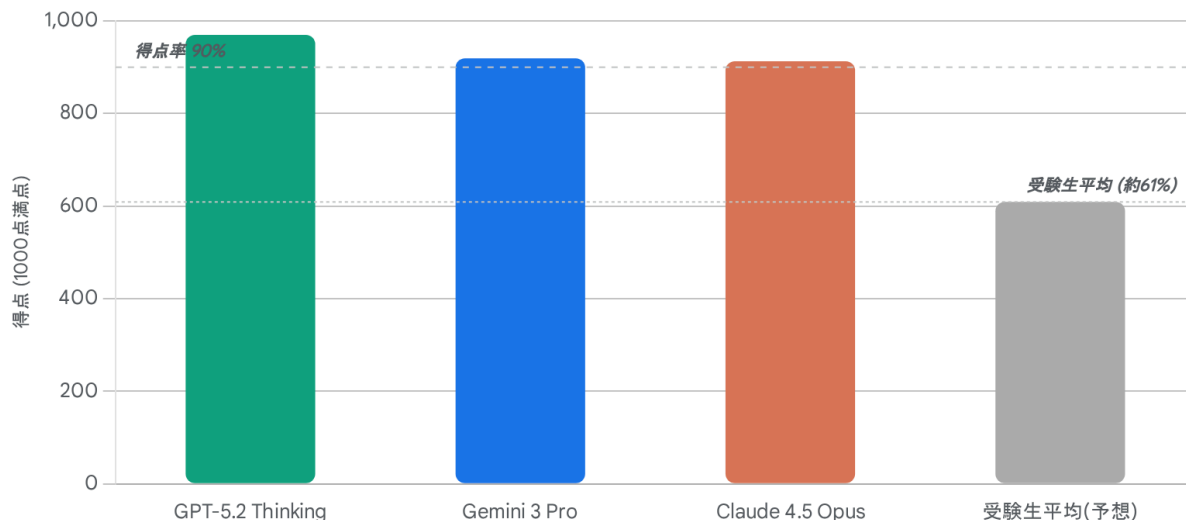
2026年1月、日本の教育界と人工知能(AI)業界は、歴史的な転換点を迎えました。AIスタートアップである株式会社LifePromptが実施した検証実験において、OpenAIの最新モデル「GPT-5.2 Thinking」が、大学入学共通テストで得点率97%という驚異的な記録を達成しました<sup>1</sup>。これは、主要15科目のうち9科目で満点(100点)を獲得し、従来のAIモデルが苦手としていた高度な論理推論や空間認識を伴う数学・理科の壁を完全に突破したことを意味します<sup>2</sup>。

過去数年間、大規模言語モデル(LLM)の進化は目覚ましく、2025年時点でも一部のモデルは難関大学合格レベルに達していました。しかし、今回の97%というスコアは、単なる「合格レベル」を超越し、人間を含む全受験者層の最上位0.1%をも凌駕する「超越的知能(Superhuman Performance)」の領域に到達したことを示唆しています。Googleの「Gemini 3 Pro」やAnthropicの「Claude 4.5 Opus」も得点率91%台を記録し、人間であれば東京大学理科三類(医学部)の足切りラインを余裕で突破する水準にありますが、GPT-5.2 Thinkingとの間には明確な性能差、いわば「6ポイントの絶対的な壁」が存在することが明らかになりました<sup>2</sup>。

本レポートでは、この衝撃的な検証結果を出発点として、各AIモデルの技術的特性、特にOpenAIが実装した「Thinking(思考)」プロセスの優位性と課題、GoogleやAnthropicのアプローチとの差異を詳細に分析します。また、共通テストという特定のベンチマークを超えて、これらのAIモデルが産業界や労働市場に与える影響、さらには「知識の記憶と再生」を主軸としてきた日本の教育システムが直面する存亡の危機についても、包括的な考察を行います。

# 2026年 大学入学共通テスト：AIモデル別総合得点比較

● GPT-5.2 Thinking ● Gemini 3 Pro ● Claude 4.5 Opus ● 受験生平均(予想)



GPT-5.2 Thinkingは97%（970点相当）を記録し、Gemini 3 ProやClaude 4.5 Opus（91%前後）を引き離しました。人間の予想平均点（約60%）と比較すると、AIの能力が完全に「受験」という枠組みを超越していることが分かります。

Data sources: [Denfaminicogamer](#), [ReseEd](#), [Ameblo \(Toshin\)](#), [Resemom](#)

## 2. 検証実験の全容と結果の多角的分析

### 2.1 圧倒的なスコアとモデル間の階層構造

LifePrompt社による2026年の検証実験は、AIモデルの性能評価における新たな基準を打ち立てました。実験結果を詳細に分析すると、現在のAIモデルは明確な階層構造を形成していることが分かります。

第一に、「超越的知能」階層に位置するのがGPT-5.2 Thinkingです。文系科目で970点（1000点満点）、理系科目で968点というスコアは、統計的な誤差の範囲を超えて、他のモデルを凌駕しています<sup>2</sup>。特筆すべきは、主要15科目のうち9科目で満点を獲得したという事実です。これには、これまでAIが苦手としてきた数学IA、数学IIBC、物理基礎、化学、情報Iなどが含まれます。

第二に、「超難関大合格レベル」階層に位置するのが、GoogleのGemini 3 Pro（919点相当）とAnthropicのClaude 4.5 Opus（913点相当）です<sup>2</sup>。これらのモデルも、人間であれば日本国内のあらゆる大学の合格圏内に入る極めて優秀な成績を収めています。しかし、GPT-5.2が完答した数学や理科の難問において、わずかな論理の飛躍や計算プロセスの不整合により失点しており、900点台前半というスコアに留まりました。

この結果は、2024年から2025年にかけて続いた「モデル性能の拮抗状態」が終わりを告げ、推論特化型アーキテクチャを採用したOpenAIが再びリードを広げたことを示唆しています。

## 2.2 満点科目の内訳と技術的ブレイクスルーの深層

GPT-5.2 Thinkingが満点を獲得した9科目のリストは、AI技術の進化の方向性を如実に示しています。

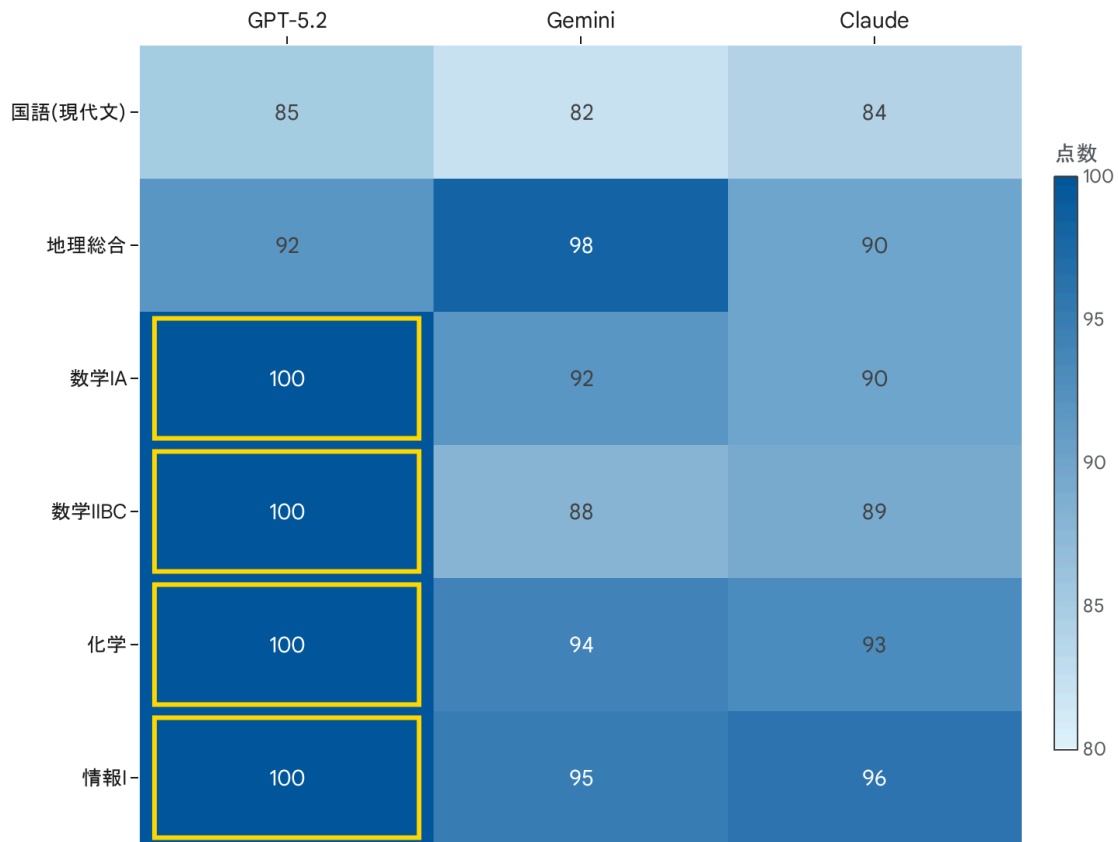
1. 数学IA
2. 数学IIBC
3. 公共、政治・経済
4. 化学
5. 物理基礎
6. 化学基礎
7. 地学基礎
8. 生物基礎
9. 情報I

特に注目すべきは、数学と理科系科目における完全制覇です。従来のLLMは、確率論的な単語予測モデルであるという性質上、厳密な論理展開や多段階の計算を要する問題において、高い確率で「もっともらしい誤答(ハルシネーション)」を生成する傾向がありました。しかし、GPT-5.2 Thinkingは、図形問題を単なるピクセルデータの集合(画像)として処理するのではなく、内部的に\*\*「座標データ」や「幾何学的制約条件」として再構築(Reconstruction)\*\*し、記号的な論理処理を行うことで正解を導き出しています<sup>2</sup>。

これは、AIの処理プロセスが、直感的なパターンマッチングを行う「システム1(速い思考)」から、熟慮と自己検証を行う「システム2(遅い思考)」へと移行したことを意味します。数学の問題を解く際、モデルは即座に回答を出力するのではなく、数分間の「思考時間」を確保し、その間に複数の解法ルートを探索・検証し、誤りを自己修正しています。このプロセスこそが、数学や物理における満点の原動力であり、GPT-5.2が「Thinking」の名を冠する所以です。

一方、情報IIにおける満点も産業的な意義が大きいと言えます。プログラミングやデータ構造、アルゴリズムに関する深い理解と、それを具体的な問題解決に適用する能力は、現代のソフトウェアエンジニアリングに直結するスキルです。AIがこの科目で満点を取れるということは、初級から中級レベルのコーディングやシステム設計タスクにおいて、AIが人間のエンジニアを代替、あるいは強力に支援できる段階にあることを裏付けています。

## AIモデル別 科目別得点ヒートマップ



GPT-5.2 Thinkingは数学・理科・情報で濃い色（満点）が並びます。一方、国語や英語リスニングでは各モデルとも色が薄くなり、AI共通の弱点であることが示唆されています。

Data sources: [Denfaminicogamer](#), [ReseMom](#)

### 2.3 Gemini 3 Proの特異点：マルチモーダル地理推論の勝利

総合点ではGPT-5.2の後塵を拝したものの、GoogleのGemini 3 Proには特筆すべき強みが見られました。それは\*\*「地理総合」における特定の視覚的推論問題\*\*です。

具体的には、南米の気候グラフと地図上の位置情報を関連付ける第3問において、Gemini 3 Proのみが正解に到達しました。他のモデルが失敗した中、Geminiは地図上の視覚情報から「アンデス山脈」という地形的特徴を正確に識別し、その標高や地形が気候に与える影響（ケッペンの気候区分など）を論理的に推論して、適切な雨温図とリンクさせることに成功しました<sup>2</sup>。

この事例は、Googleが長年蓄積してきたGoogle MapsやGoogle Earthなどの地理空間データと、検索エンジン由来の膨大な知識ベースが、Geminiのマルチモーダル学習に深く統合されていること

を示唆しています。視覚情報(地図画像)と言語情報(気候の説明)を高度に融合させる能力において、特定のドメインではGeminiがGPT-5.2を上回る可能性があることは、今後のモデル選定や用途開発において重要な知見となります。

### 3. 残された課題:「感情」と「視覚的論理」の壁

得点率97%という数字は、AIがもはや大半の知的タスクにおいて人間を凌駕したことを示していますが、残りの3%の失点領域にこそ、AIの現状の限界と、人間知能の独自性が凝縮されています。LifePrompt社の分析により、全モデルに共通する2つの主要な弱点が浮き彫りになりました。

#### 3.1 国語(小説)における「感情の機微」の理解不能

第一の壁は、国語の小説問題です。論説文や実用的な文章の読解において、AIは人間以上の要約力と論理構成の把握能力を発揮し、ほぼ満点を獲得しました。しかし、小説における「登場人物の心情の変化」や「文脈に埋め込まれた感情の機微(subtleties of emotion)」を問う問題では、全てのモデルが苦戦を強いられました<sup>2</sup>。

例えば、登場人物の何気ない発言や行動の裏にある葛藤、皮肉、あるいは愛情といった複雑な心理状態を、文脈から読み取る能力です。AIは膨大なテキストデータから「悲しいときは泣く」「怒っているときは叫ぶ」といった一般的なパターンを学習していますが、小説独特の「言外のニュアンス」や、個人の経験や文化的背景に依存する共感的な理解までは到達していません。GPT-5.2であっても、論理的に感情を推測しようと試みますが、人間の持つ身体性や実存的な経験に基づく直感的な共感とは異なり、微妙なニュアンスを取り違えるケースが散見されました。これは、AIが「情報(Information)」を処理できても、「情動(Affect)」を主観的に体験していないという根本的な制約に起因する、いわゆる「記号接地問題(Symbol Grounding Problem)」の一形態とも言えるでしょう。

#### 3.2 英語リスニングにおける「視覚的論理」の断絶:バス問題の衝撃

第二の壁は、英語リスニング問題における視覚情報処理の失敗です。特に象徴的だったのが、バスの乗降方法に関する問題です<sup>2</sup>。

この問題では、音声スクリプト(実験ではテキストとして入力)において、「後ろのドアから乗って、前のドアから降りる(Enter from the rear, exit from the front)」というルールが明確に述べられていました。驚くべきことに、GPT-5.2を含むすべてのモデルの「思考ログ(Thinking Log)」を確認すると、AIはこのルールをテキストとして正しく認識し、文字起こしまで行っていました。論理的な理解という点では、AIは正解に到達していたのです。

しかし、最終的な回答として、このルールに対応するイラスト(バスのドアと矢印が描かれた図)を選択する段になると、全てのモデルが誤ったイラストを選んでしまいました。これは、「言語的な論理理解」と「視覚的な表現へのマッピング(Grounding)」の間に深刻な断絶があることを示しています。AIは「後ろから乗る」という概念を言語的には理解していても、それを「平面図上のバスの形状、ドアの位置、矢印の方向」という視覚的シンボルと正確に結びつける空間的・視覚的推論能力において、致命的なエラーを起こしました。

Gemini 3 Proのようにマルチモーダル性能を売りにするモデルであっても、このような抽象的な図解

やダイアグラムの解釈において躓く事実は、自動運転やロボティクスなど、現実世界の視覚情報を論理的に処理する必要がある分野への応用において、依然として解決すべき課題が残されていることを警告しています。

## 4. アーキテクチャの進化とコスト・パフォーマンスの相関

### 4.1 「Thinking」プロセスの代償：時間とコスト

GPT-5.2が圧倒的なスコアを叩き出した背景には、OpenAIが導入した新たな推論パラダイムである「Thinking」プロセスがあります。従来のモデルが入力に対して即座に確率的なトークン予測を行っていたのに対し、GPT-5.2 Thinkingは、回答を出力する前に内部的な思考プロセス(Chain of Thought)を実行します。

このプロセスにより、モデルは自身の生成した中間的な推論ステップを検証し、誤りがあれば自己修正(Self-Correction)を行うことが可能になります。数学の計算ミスや論理の飛躍が激減したのは、このメカニズムの恩恵です。

しかし、この高性能には明確な代償が伴います。それは「時間」と「計算リソース」です。LifePrompt社の実験データによると、\*\*GPT-5.2 Thinkingが全科目の解答に要した時間は約5時間30分(330分)でした<sup>4</sup>。これに対し、Gemini 3 ProやClaude 4.5 Opusは、わずか1時間40分(100分)\*\*程度で試験を完走しています。人間の受験生の規定試験時間が約10時間であることを考えれば、GPT-5.2でも十分に高速ですが、GeminiやClaudeと比較すると3倍以上の時間を要していることになります。

この「時間差」は、そのままAPIコストの差にも直結します。Thinkingモデルは、ユーザーに見えない「思考トークン」を大量に消費するため、推論コストが通常のモデルの数倍から十数倍に膨れ上がる傾向があります<sup>5</sup>。したがって、産業応用の観点からは、97%の精度が必要なクリティカルなタスク(医療診断の補助、契約書の法的精査、複雑なシステム設計など)にはGPT-5.2 Thinkingが適していますが、速度とコスト効率が重視される一般的なタスク(要約、定型メール作成、単純なデータ処理)においては、Gemini 3 ProやClaude 4.5 Opusの方が「コストパフォーマンスの良い選択肢」となるでしょう。この「適材適所」の使い分けが、2026年のAI活用の鍵となります。

### 4.2 Gemini 3 Proの課題：熱暴走とコンテキスト保持の脆弱性

一方、GoogleのGemini 3 Proに関しては、別の技術的課題が浮き彫りになっています。それはサーマルスロットリング(熱暴走)に起因する性能低下です。一部のユーザー報告によると、モバイルデバイスや冷却能力の低いエッジ環境において、Gemini 3 Proを長時間稼働させるとデバイス温度が急上昇し、それに伴って推論能力が著しく低下する現象が確認されています<sup>6</sup>。

具体的には、長文のコンテキスト(文脈)を保持し続けることが困難になり、会話の初期に提示された情報を忘却したり、論理的な整合性を欠く回答を生成したりする「コンテキスト脱落」が発生しやすくなります。クラウドAPI経由でのベンチマークテストでは、サーバー側の強力な冷却環境によりこの問題は顕在化しにくいですが、スマートフォンやPCなどのローカル環境で動作する「オンデバイスAI」としての展開を視野に入れた場合、この熱問題は致命的な欠陥となり得ます。Googleが今後、モデルの軽量化や推論効率の最適化(蒸留モデルの開発など)によってこの課題をどう克服するかが、



Gemini普及の試金石となるでしょう。

## 5. 東京大学「理科三類」合格ラインへの到達と社会的衝撃

### 5.1「足切り」概念の崩壊と入試システムの形骸化

今回の検証結果が教育界に与える最大の衝撃は、日本最難関とされる東京大学理科三類(医学部)の第一次選抜(足切り)ラインを、AIが完全に突破したという事実です。

東大理科三類の足切りラインは年度によって変動しますが、例年おおむね得点率\*\*77%~80%程度で推移しています。合格者の平均得点率を見ても88%~90%程度です<sup>7</sup>。これに対し、GPT-5.2 Thinkingの97%はもちろんのこと、**Gemini 3 Pro**や**Claude 4.5 Opus**の91%\*\*台というスコアも、このラインを余裕で上回っています。

これは何を意味するのでしょうか。それは、マークシート方式で基礎学力を測る「大学入学共通テスト」という試験制度において、AIが人間(国内トップ層の受験生を含む)を完全に凌駕したということです。かつて「AIには東大合格は無理だ」と言われた時代がありましたが、2026年現在、少なくとも一次試験の段階では、AIは「優秀な受験生」ではなく「超人的な受験生」としての地位を確立しました。この事実は、知識の正確な再生や定型的な論理処理を測るテストが、もはや人間の知的能力の測定尺度として機能しなくなりつつあることを示唆しています。

### 5.2「東大二次試験」への挑戦と記述能力の進化

LifePrompt社は、共通テストでの圧勝を受けて、次なるステップとして「AI vs 東大二次試験」プロジェクトを進行させています<sup>3</sup>。共通テストが選択式であるのに対し、二次試験は高度な記述力が求められる論述式です。

これまでのAIは、長文の論述において論理が一貫しなかったり、数学の証明問題で途中経過を省略しすぎたりすることで減点される傾向がありました。しかし、GPT-5.2 Thinkingが数学で見せた「論理の再構築」能力や、自己検証プロセスによる厳密性の向上は、記述式試験においても満点に近い解答を作成できる可能性を強く示唆しています。もしAIが東大の二次試験、特に数学や理科の難問記述において合格点を叩き出すようになれば、「思考力や表現力を問う記述式ならばAIに対抗できる」という人間の最後の砦さえも崩れ去ることになります。

## 6. 教育への示唆: 私たちは何を学ぶべきか?

### 6.1「知識」の価値暴落と「検証力(Audit)」の高騰

共通テストで97%を取れるAIが、月額20ドル(約3,000円)程度で誰でも利用可能になる社会において、人間が「正解を知っていること」の経済的価値は暴落します。百科事典的な知識を暗記し、それを正確にアウトプットする能力は、もはやAIの独壇場であり、人間がそこで勝負することは非効率極まりない行為となります。

その一方で、極めて高い価値を持つようになるのが\*\*「検証力(Audit Capability)」\*\*です。AIは97%の正解を出しますが、残りの3%で、先述のバス問題のような「人間なら犯さないような奇妙なミス」を犯す可能性があります。また、もっともらしい嘘(ハルシネーション)をつくりリスクもゼロではありません。

ん。したがって、AIが出力した高度な回答や成果物に対して、「本当に正しいのか？」「前提条件に誤りはないか？」「倫理的な問題はないか？」を批判的に精査し、最終的な責任を持って判断を下す能力が、これからの人間に求められる最も重要なスキルとなります。教育の現場では、正解を出す訓練よりも、AIの回答を疑い、検証し、修正する訓練（AIオーディティング教育）へのシフトが急務です。

## 6.2「思考体力」の重要性和AI協働

GPT-5.2が時間をかけて「Thinking」を行うことで高得点を得た事実は、人間にとっても重要な教訓を含んでいます。即座に直感的な答えを出す「瞬発力（System 1）」の領域では、人間はもはやAIの速度と精度に勝てません。しかし、複雑な問題に対して粘り強く向き合い、多角的な視点から仮説を立て、時間をかけて論理を積み上げる\*\*「思考体力（System 2）」\*\*の重要性は、むしろ高まっています。

AIは強力なツールですが、問いを立てるのは依然として人間の役割です。「何を解くべきか」という課題設定能力、そしてAIという異質の知能と対話しながら、自身の思考を拡張していく「AI協働能力」こそが、2026年以降の社会で生き残るための必須条件となるでしょう。新井紀子氏らが以前から指摘していたように、「AIにはできない読解力」や「意味の真の理解」を問う教育への回帰は、もはや理想論ではなく、生存戦略としての必然性を帯びています<sup>10</sup>。

## 7. 結論：人間は「回答者」から「監督者」へ

2026年の大学入学共通テストにおけるAIの圧勝は、単なる技術的なマイルストーン以上の意味を持っています。それは、近代教育システムが前提としてきた「知識の習得と再生」という能力評価モデルの終焉を告げる警鐘です。

- **GPT-5.2 Thinking**は、数学・理科・情報を「解決済み（Solved）」の課題とし、高度な推論能力を民主化しました。
- **Gemini 3 Pro**は、視覚と知識の融合という新たな可能性を示しつつ、ハードウェア制約という現実的な課題も浮き彫りにしました。
- しかし、感情の機微や実世界での視覚的論理という、人間にとって当たり前の領域には、依然としてAIが越えられない壁が存在します。

これからの時代、人間はテストの「回答者（Answerer）」としての役割をAIに譲り渡し、AIという強力な知能リソースを指揮・監督し、その出力の価値とリスクを判断する\*\*「監督者（Director）」\*\*としての役割を担うことになります。2026年の共通テストは、その役割転換が不可逆的に始まったことを告げる、歴史的な特異点として記憶されることになるでしょう。教育機関、企業、そして私たち個人は、この新しい現実に適応するために、学習と評価のシステムを根本から再設計する必要があります。

## 参考文献データソース

本レポートの分析は、以下の検証データおよび報道記事に基づいています。

- 検証実施: 株式会社LifePrompt (2026年1月20日公表)
- 主要情報源:
  - 日本経済新聞:<sup>1</sup>
  - Denfaminicogamer:<sup>2</sup>



- Resemom: <sup>3</sup>
- LLM-Stats / Benchmarks: <sup>12</sup>
- その他予備校データ・専門家コメント: <sup>5</sup>

## 引用文献

1. 大学入学共通テスト、OpenAIは9科目満点 得点率97%でGoogleに勝利, 1月 20, 2026 にアクセス、  
<https://b.hatena.ne.jp/entry/s/www.nikkei.com/article/DGXZQOUC190LP0Z10C26A1000000/>
2. 2026年の大学入学共通テスト、ChatGPTが「9科目で満点」を獲得。合計点数でも Gemini、Claudeに差をつける, 1月 20, 2026にアクセス、  
<https://news.denfaminicogamer.jp/news/260120g>
3. 共通テスト、AIが9科目で満点...図形や濃淡に課題 | 教育業界 ..., 1月 20, 2026にアクセス、  
<https://reseed.resemom.jp/article/2026/01/20/12501.html>
4. 【共通テスト2026】AIが9科目で満点...図形や濃淡に課題 - リセマム, 1月 20, 2026にアクセス、  
<https://resemom.jp/article/2026/01/20/84722.html>
5. GPT-5.2 vs Gemini 3 Pro: Which AI Model is Better in 2026 ..., 1月 20, 2026にアクセス、  
<https://evolink.ai/blog/gpt-5-2-vs-gemini-3-pro-comparison-2026>
6. Gemini 3 significantly worse than 2.5 Pro at long context ..., 1月 20, 2026にアクセス、  
<https://discuss.ai.google.dev/t/gemini-3-significantly-worse-than-2-5-pro-at-long-context-temperature-likely-to-blame/110888>
7. 2026東京大学理科三類の共通テストボーダー・足切り、合格者最低点, 1月 20, 2026にアクセス、  
<https://igakubuyobiko.com/bible/?p=6796>
8. 【2026年度版】東京大学 理科Ⅲ類(医学部) 共通テスト足切り(第 ..., 1月 20, 2026にアクセス、  
<https://shotscha.com/medical-school/information-by-university/score/u-tokyo-3>
9. MEDIA - ライフプロンプト, 1月 20, 2026にアクセス、  
<https://lifeprompt.net/news/category/media>
10. ニュース - 教育のための科学研究所 - リーディングスキルテスト, 1月 20, 2026にアクセス、  
[https://rst-web.s4e.jp/blogs/blog\\_entries/tag/17/id:8/limit:100?frame\\_id=17](https://rst-web.s4e.jp/blogs/blog_entries/tag/17/id:8/limit:100?frame_id=17)
11. AIに読解力があると思う人に知ってほしい現実 学生の新常識は「シン ..., 1月 20, 2026 にアクセス、  
<https://toyokeizai.net/articles/-/370228?display=b>
12. LLM Benchmarks 2026 - Complete Evaluation Suite, 1月 20, 2026にアクセス、  
<https://llm-stats.com/benchmarks>
13. AI Leaderboards 2026 - Compare LLM, TTS, STT, Video, Image ..., 1月 20, 2026にアクセス、  
<https://llm-stats.com/>
14. 【共通テスト2026】文系予想平均点、予備校間の乖離に絶句, 1月 20, 2026にアクセス、  
<https://ameblo.jp/doushitaraii/entry-12953908315.html>