

# 次世代推論モデル「Gemini 3.1 Pro」の全容とAGIへの到達点：競合フロンティアモデル比較およびDeep Thinkとの構造的差異

## Gemini 3.1 pro

### 1. 序論：推論時間計算量スケーリングの時代へのパラダイムシフト

2026年2月は、人工知能の進化史、とりわけ大規模言語モデル(LLM)の発展において極めて重要な転換点として記録される。Googleによって発表された「Gemini 3.1 Pro」、およびそれに先行して導入された「Gemini 3 Deep Think」の登場は、AIの性能向上が「パラメータサイズの拡大(学習時の計算量)」から「推論時間計算量の拡大(Inference-Time Compute Scaling)」へと完全にシフトしたことを決定づけた<sup>1</sup>。

従来のモデルアップデートが「.5」という刻み幅で行われてきたのに対し、今回Googleが「.1」という異例のバージョン体系を採用したことは、本モデルが機能の単純な横展開(新たなモダリティの追加など)ではなく、中核となる「推論(Reasoning)エンジン」の知能底上げに特化していることを示唆している<sup>2</sup>。現代の科学、研究、ソフトウェアエンジニアリングにおける課題は、単一のプロンプトに対する単純な回答では解決できない水準に達しており、複数の前提条件を統合し、長期間にわたる計画を立て、自律的に検証を繰り返す能力が求められている<sup>2</sup>。

本報告書は、Gemini 3.1 Proのアーキテクチャ特性、AGI(汎用人工知能)への登竜門とされる「ARC-AGI-2」ベンチマークにおける歴史的成果、Claude Opus 4.6やGPT-5.3 Codex等の最先端競合モデルとの詳細な比較分析、そして最上位の特殊推論モードである「Deep Think」との構造的・戦略的差異について、包括的かつ深層的な分析を提供する。

### 2. Gemini 3.1 Proのアーキテクチャと基本性能

Gemini 3.1 Proは、複雑な多段階の課題解決に特化して設計されたネイティブマルチモーダル・フロンティアモデルである。前モデルであるGemini 3.0 Proと同一の価格設定を維持しながら、その実質的な能力は世代交代レベルの飛躍を遂げている<sup>3</sup>。

#### 2.1. コア仕様とエージェント向け最適化エンドポイント

本モデルは、最大100万トークンの入力コンテキストウィンドウと、明示的に拡大された65,000トークンの出力上限を備えている<sup>7</sup>。入力データ型としてはテキスト、画像、動画、音声に加えてPDFをネイティブにサポートしており、新たに最大100MBまでのファイルアップロードが可能となったことで、巨大なコードベースや学術論文群の一括処理能力が大幅に向上した<sup>7</sup>。また、YouTubeのURLを直接プ

プロンプトに記述することで、事前のダウンロードなしに動画コンテンツのコンテキストを解析する機能も実装されている<sup>8</sup>。

価格体系は、100万入力トークンあたり2.00ドル(200Kトークン未満の場合)、100万出力トークンあたり12.00ドルと、前世代から据え置かれている<sup>8</sup>。しかし特筆すべきは、出力効率の劇的な改善である。Gemini 3.1 Proは、前モデルと比較して出力トークン消費量を平均して約10%から15%削減しつつ、より信頼性の高い結果を生成する<sup>8</sup>。JetBrains社のAIディレクターによる実世界の評価でも、この効率化により、同じタスクを実行する際の実際のトークン消費量が減少し、企業規模でのAPI運用において大幅なコスト削減をもたらすことが確認されている<sup>8</sup>。これは、モデルが単に「短い回答を返す」ようになったのではなく、内部の推論プロセスが最適化され、冗長な説明を省いて核心を突くコードやロジックを直接出力する「思考の洗練」が達成されたことを意味する。

さらに、マルチステップのワークフローや自律型エージェント開発向けに、新たにgemini-3.1-pro-preview-customtoolsという専用エンドポイントが提供されている<sup>7</sup>。このエンドポイントは、view\_fileやsearch\_codeといったカスタムツールの呼び出し優先度を最適化しており、Bashコマンドと独自関数を頻繁に行き来する複雑な自律型タスクにおいて、汎用モデル特有のツール呼び出しの不安定さを解消し、高い完遂能力を発揮する<sup>7</sup>。

## 2.2. 知能の応用: クリエイティブコーディングとシステム統合の深淵

Gemini 3.1 Proの推論能力は、抽象的なロジックを具体的な成果物に変換するプロセス、すなわち「Vibe Coding(意図や雰囲気を読み取ったコーディング)」において顕著に表れる<sup>2</sup>。単なる構文規則の出力ではなく、ユーザーの意図、製品のスタイル、そしてドメイン特有のコンテキストを深く理解した上でのコード生成が可能となっている<sup>3</sup>。

具体的な応用例として、コードベースのアニメーション生成が挙げられる。テキストプロンプトから直接、Webサイトで即座に利用可能なアニメーションSVGを生成する機能である<sup>2</sup>。ピクセルベースの動画ファイル(MP4やGIFなど)とは異なり、純粋なコードで構築されるため、いかなる解像度への拡大・縮小にも劣化せず、ファイルサイズも極めて軽量に保たれる<sup>2</sup>。変形する3Dの段ボール箱や、詳細な動作指定を含むアニメ風キャラクターの生成など、空間的・数学的な計算が不可欠なタスクを単一のプロンプトで完遂する能力は、従来のモデルには見られなかった特質である<sup>3</sup>。

また、複雑なデータストリームの統合と可視化においても高度な能力を示す。例えば、公開されている国際宇宙ステーション(ISS)のテレメトリAPIを自律的に解析・設定し、リアルタイムの軌道を可視化する高精度な3Dダッシュボードを構築することができる<sup>5</sup>。このプロセスには、現在のUTC時刻に基づいた太陽の位置計算(地球の昼夜の正確なレンダリング)といった複雑なロジックが組み込まれており、APIの仕様理解と物理的・数学的推論の高度な統合を証明している<sup>5</sup>。

さらに、文学的テーマのインターフェース翻訳という、極めて抽象度の高いタスクも実行可能である。エミリー・ブロンテの小説『嵐が丘』のトーンや雰囲気を推論し、その文学的本質を捉えた現代的なパーソナルポートフォリオ(Webサイト)のUI/UXを設計・制作する能力は、本モデルが単なる技術的要件の処理を超え、深い意味的・美学的な理解を備えていることを示している<sup>5</sup>。その他にも、ハンドトラッキングなどのユーザー操作に連動し、動きに合わせて変化する生成音楽を組み合わせた「ムク

「ドリ」の群舞」の没入型3Dシミュレーションの構築など、研究者やデザイナーのプロトタイピングプロセスを根本から変革するポテンシャルを秘めている<sup>5</sup>。

### 3. AGIへの登竜門:「ARC-AGI-2」における歴史的飛躍

LLMの性能評価において、MMLUなどの従来のベンチマークはすでに飽和状態にあり、インターネット上のデータセットを丸暗記すること(データコンタミネーション)によるスコアのインフレが業界の深刻な課題となっていた<sup>15</sup>。この限界を打破し、真の「流動性知能(Fluid Intelligence)」を測定するためにフランソワ・ショレ(François Chollet)らによって提唱されたのが「ARC-AGI」である<sup>16</sup>。

#### 3.1. ARC-AGI-2の過酷さと「流動性知能」の定義

2025年に導入された「ARC-AGI-2」は、AIモデルが事前に学習したパターンマッチングや力技(ブルートフォース)で解くことが不可能なように設計された、極めて厳格なベンチマークである<sup>16</sup>。少数の訓練例(色付きのグリッドの入力と出力のペア)から背後にある「未知の変換ルール」を抽象化し、全く新しいテストケースに適用する能力が問われる<sup>17</sup>。人間のテスト参加者は、専門的なプログラミング知識や高度な数学的背景がなくても、平均して100%近い正答率(あるいは1タスクあたり数分で解ける水準)を示す一方で、純粋なLLMアーキテクチャは長らく1桁台のスコアに留まっていた<sup>18</sup>。

この人間とAIの圧倒的な乖離は、現在のAIが「記号に意味を割り当てる(Symbolic Interpretation)」能力や、「複数の相互作用するルールを同時に適用する(Compositional Reasoning)」能力に致命的な弱点を持っていたことに起因する<sup>19</sup>。AIは画像の対称性や反転をチェックすることはできても、記号そのものに意味論的な重要性を見出すことができず、結果として未知の論理体系への適応に失敗していたのである<sup>19</sup>。

#### 3.2. スコア77.1%が意味するブレイクスルーと業界への衝撃

この過酷な条件下において、Gemini 3.1 Proは「ARC-AGI-2」で検証済みスコア77.1%という驚異的な記録を達成した<sup>2</sup>。前モデルであるGemini 3.0 Proのスコアが31.1%であったことを考慮すると、わずかな期間で推論能力が2.4倍以上に跳ね上がったことになる<sup>6</sup>。

評価モデル	ARC-AGI-2 スコア	前モデル(Gemini 3.0 Pro)からの成長率
Gemini 3.1 Pro	77.1%	+148%
Claude Opus 4.6	68.8%	N/A
GPT-5.2	52.9%	N/A
Gemini 3.0 Pro	31.1%	基準値

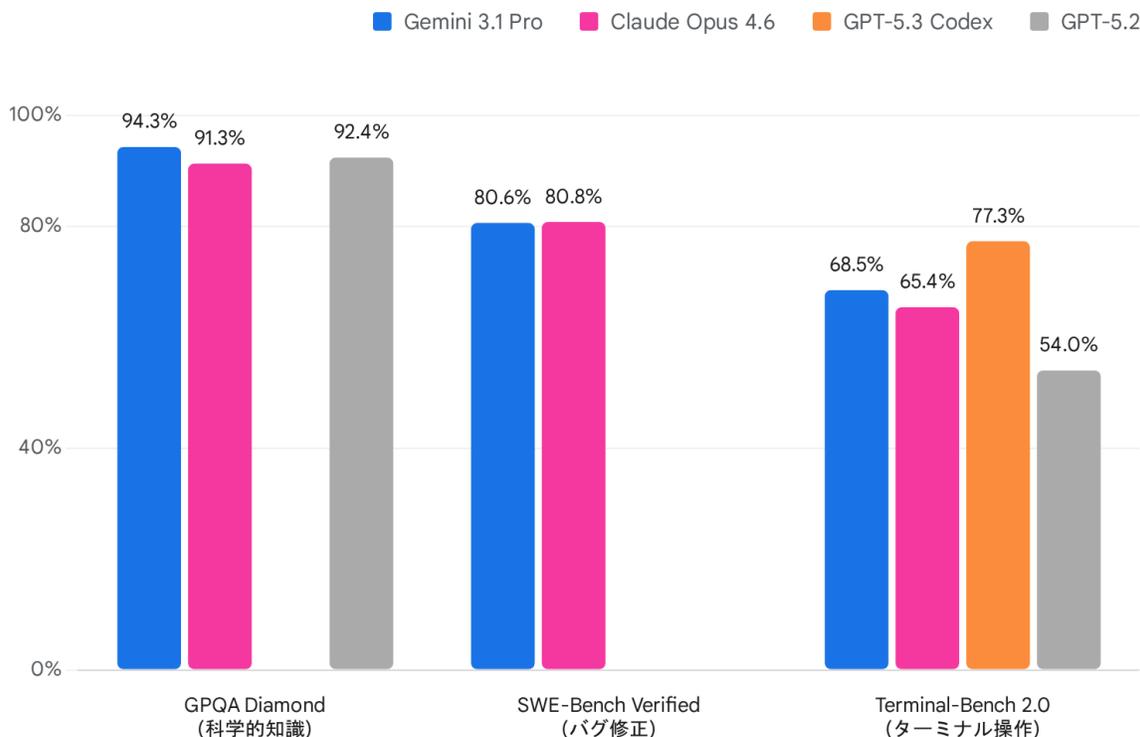
この結果は、AIの発展史において極めて重要な意味を持つ。「単にモデルのパラメータサイズを対数的に拡大してもARC-AGI-2を解くことはできない(Scale is Not Enough)」という業界の定説に対し、Gemini 3.1 Proは推論レイヤーの根本的な最適化によって厚い壁を突破したのである<sup>19</sup>。

もちろん、77.1%という数字はまだ100%には達しておらず、ARC-AGI-2が「完全に解決された」ことを意味するものではない<sup>16</sup>。ベンチマークの開発者たちが指摘するように、真の意味でのAGIと呼ぶためには、人間と同等の100%に近い成功率が求められるからである<sup>16</sup>。しかし、競合モデルであるClaude Opus 4.6(68.8%)やGPT-5.2(52.9%)を明確に凌駕し、AIが未知のタスクに対して柔軟に適応する「適応性(Adaptability)」を大規模に獲得し始めた事実は、2026年内のAGIベンチマーク飽和を予感させるに十分なマイルストーンであると言える<sup>1</sup>。

## 4. 競合4モデルの比較分析: 強みと弱みの解剖

現在のフロンティアAI市場は、Google、Anthropic、OpenAIによる熾烈な三つ巴の様相を呈している。本セクションでは、Gemini 3.1 Proを、その直接の競合である「Claude Opus 4.6」、「GPT-5.3 Codex」、および「GPT-5.2 xhigh」と多角的に比較し、各モデルの特質と最適なユースケースを浮き彫りにする<sup>1</sup>。

# 主要フロンティアモデルのドメイン別性能比較



Gemini 3.1 Proは科学的知識（GPQA）でトップに立つが、実世界のバグ修正（SWE-Bench）ではClaude Opus 4.6と拮抗し、ターミナル操作（Terminal-Bench）ではGPT-5.3 Codexに遅れをとっている。

データソース: [ChatGPT Lab](#), [Digital Applied](#), [DeepMind](#), [NXCode](#)

## 4.1. 推論の「広さ(Breadth)」と「深さ(Depth)」の対立構造

Gemini 3.1 ProとClaude Opus 4.6の比較は、AIアーキテクチャ開発における設計思想の根源的な違いを鮮明に示している。

Gemini 3.1 Proは、「広範な推論(Breadth)」と「複数ツールの同時並行的な調整」において圧倒的な強みを持つ<sup>1</sup>。アルゴリズム問題解決能力を測るLiveCodeBench Proにおいて、Gemini 3.1 Proは2887 Eloという歴代最高水準のスコアを記録し、他の全モデルを引き離れた<sup>1</sup>。また、PhD級の科学的推論を問うGPQA Diamondでは94.3%という驚異的な新記録を樹立している<sup>1</sup>。さらに、エージェントが複数のツールを同時にオーケストレーションする能力を測るMCP Atlasでは69.2% (Opusは59.5%)、完全自律型のマルチステップ・タスク実行を測るAPEX-Agentsでは33.5% (Opusは29.8%)、自律的なウェブ検索とPython実行を組み合わせたBrowseCompでは85.9%を記録し、複雑なワークフ

ロー全体を俯瞰する能力の高さを示した<sup>1</sup>。

一方で、Claude Opus 4.6は「深さ(Depth)」と特定のドメインにおける「厳密な専門性」において依然として王座を保っている<sup>1</sup>。実世界のオープンソースリポジトリのバグを自律的に修正する能力を測るSWE-Bench Verifiedにおいて、Gemini 3.1 Proは80.6%と極めて高い水準に達したが、Opus 4.6の80.8%には僅か0.2ポイント及ばなかった<sup>1</sup>。また、高度なドメイン知識と専門家の推論プロセスを模倣する専門的タスク(GDPval-AA)においては、Opus 4.6が1606 Eloを記録し、Gemini 3.1 Pro(1317 Elo)を圧倒している<sup>1</sup>。この結果は、実稼働環境でのミッションクリティカルなバグ修正や、深いドメイン知識が必要な専門的リサーチにおいては、Opusのアーキテクチャが依然として高い信頼性を持つことを示唆している。

同一プロンプトを用いたクリエイティブタスクの比較検証(VibeCheck)でも、この傾向は裏付けられている。和モダン茶室のランディングページデザインやペリカンのSVGアニメーション生成といったビジュアルと空間認識が絡むタスクでは、Gemini 3.1 Proは洗練されたUIや滑らかな動きを実現し、Opusに匹敵または凌駕する能力を見せた<sup>7</sup>。しかし、レトロな3D宇宙船ゲームの生成といった、極めて複雑なロジックが絡み合う安定性が求められるタスクにおいては、Gemini 3.1 Proはエラーを引き起こし出力を完遂できないケースが報告されており、ビジュアル生成の強さと純粋なロジック安定性との間にギャップが存在することが指摘されている<sup>7</sup>。

## 4.2. エージェント機能とターミナル操作に残る課題

自律型エージェントのパフォーマンスにおいて、Gemini 3.1 Proは劇的な進化を遂げたものの、オペレーティングシステムレベルの操作においては特定のモデルの後塵を拝している。ターミナル環境を通じたエージェントのコーディング能力(環境構築、シェルスクリプトの実行、システムレベルのデバッグなど)を測定するTerminal-Bench 2.0において、Gemini 3.1 Proのスコアは68.5%であった。これは前モデルの56.9%からは20%以上の相対的な大幅改善であるものの、この領域に特化してチューニングされたOpenAIのGPT-5.3 Codexは77.3%を叩き出しており、約9ポイントの決定的な乖離が存在する<sup>3</sup>。OSレベルでの深い相互作用や、純粋なコマンドラインベースのデバッグ環境においては、Codexアーキテクチャの堅牢性が揺らいでいないのが現状である<sup>3</sup>。

## 4.3. 超長文脈における情報抽出精度の低下(Lost in the Middle)

Gemini 3.1 Proのアーキテクチャ上の最大の盲点として指摘すべきは、100万トークン規模の超長文脈(Long-context)入力に対する精緻な情報抽出精度の低下である<sup>7</sup>。膨大なコンテキストの中から特定の情報を探し出す能力を測るMRCR v2(Needle-in-a-haystack)テストにおいて、128Kトークンまでは84.9%という高い精度を維持するものの、入力が100万トークン規模に達すると、ピンポイントでの情報検索精度が26.3%にまで急落することが検証データから明らかになっている<sup>7</sup>。

これは、大容量のコンテキストウィンドウが物理的に「入力可能」であることと、その全体から「正確に情報を引き出せる」ことが同義ではないことを厳格に示している。巨大なモノリスリポジトリの全コードベースや、数千ページに及ぶ法的文書を一度に丸ごと入力して微細な事実関係を検索させる用途においては、モデル単体のコンテキストウィンドウに依存するのではなく、外部のベクトルデータベースを用いたRAG(検索拡張生成)システムや、適切なチャンキング技術を併用したアーキテクチャ設

計が依然として不可欠であると言える。

#### 4.4. 経済性と「モデル・ルーティング戦略」の台頭

能力面での局地的な拮抗とは対照的に、コスト構造においてはGemini 3.1 Proが破壊的な優位性を持っている。100万トークンあたりの入力コストは、Gemini 3.1 Proが2.00ドルであるのに対し、Claude Opus 4.6は15.00ドルであり、ここに実に7.5倍もの価格差が存在する<sup>1</sup>。出力コストにおいても、Gemini 3.1 Pro(12.00ドル)はOpusの価格を大きく下回っている<sup>1</sup>。

この圧倒的な経済的差異は、エンタープライズアーキテクチャにおける「インテリジェント・モデル・ルーティング」という新たなシステム設計パラダイムを必須のものとしている。大規模な自律型コーディングエージェントや、数百万トークンを日々消費する大量の文書解析ワークフローにおいて、すべての処理をOpus 4.6に任せることは経済的に非現実的である。したがって、全体の90%を占める一般的な推論、アルゴリズム設計、データ統合、マルチツールオーケストレーションタスクを安価で高速なGemini 3.1 Proで処理し、極めて高い精度や深いドメイン知識が要求される残りの10%(本番環境での致命的なバグ修正や、専門的な法務・財務の推論など)のみをClaude Opus 4.6にルーティングするというハイブリッド構成が、現在のプロダクション環境におけるベストプラクティスとして確立されつつある<sup>1</sup>。

### 5. 究極の知能「Gemini Deep Think」との構造的・戦略的差異

Gemini 3.1 Proの能力を正確に位置づけ、真のポテンシャルを引き出すためには、一足先の2026年2月12日に発表された特殊推論モード「Gemini 3 Deep Think」とのアーキテクチャ上の違いを深く理解することが不可欠である<sup>1</sup>。世間では両者が混同されがちであるが、これらは本質的に異なるアプローチに基づくシステムである。

#### 5.1. 推論時間計算量(Inference-Time Compute)のパラダイム

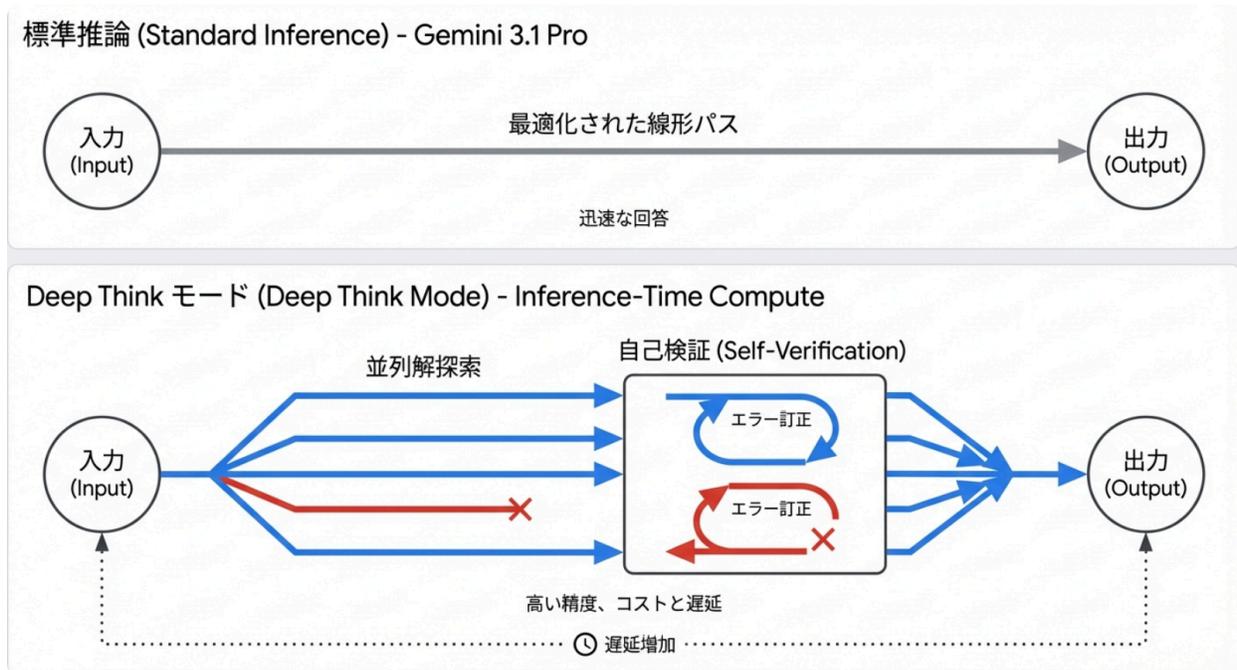
極めて重要な前提として、Gemini 3 Deep Thinkは、新しい独立した大規模なパラメータを持つ基盤モデル(重みネットワーク)ではない。標準のGemini 3アーキテクチャの上に構築された\*\*「推論レイヤー」\*\*であり、出力を生成する際の「計算資源と時間(Inference-time compute)」を意図的に増大させるアルゴリズムの総称である<sup>1</sup>。これは、モデルのサイズを大きくして知識を詰め込む従来の学習パラダイムから、推論時により長く深く考えさせることで困難な課題を解くというパラダイムへの劇的な転換である<sup>1</sup>。

通常のLLMが入力に対して一度の前向き伝播(Single forward pass)で即座に回答を生成するのに対し、Deep Thinkモードは内部で以下のような3段階の複雑なパイプラインを実行する<sup>1</sup>。

1. 問題の分解(Problem Decomposition): 複雑なクエリを受け取ると、直ちに回答を生成するのではなく、内部の推論チェーンで問題を複数のサブ問題に分割し、解決のための最適な戦略を計画する。例えば高度な数学の問題に対して、関連する定理を特定し、論理的なステップを事前に構築する<sup>1</sup>。

2. 並列解探索 (Parallel Solution Search) : 同時に複数の仮説や解決パスを並行して生成・検討する。例えば競技プログラミングにおいて、動的計画法、貪欲アルゴリズム、グラフ理論といった異なるアプローチを並列でシミュレーションし、どれが最適解にたどり着くかを探索する(モンテカルロ木探索的なアプローチ)<sup>1</sup>。
3. 検証と出力 (Verification and Output) : 選択された解決パスに対して自己検証 (Self-verification) を行い、論理的な一貫性に欠陥がないか、エッジケースで破綻しないかを確認する。エラーが発見されれば訂正ループを回し、すべての検証を通過した後にのみ、最終的な回答を出力する<sup>1</sup>。

## 推論時間計算量 (Inference-Time Compute) のアーキテクチャ構造



Gemini 3.1 Proが最適化された線形パスで迅速に回答を生成するのに対し、Deep Thinkモードは内部で複数の仮説を並列に探索・自己検証 (エラー訂正ループ) するため、高い精度と引き換えに計算コストと遅延が増加する。

この多大な計算資源を投じた推論プロセスにより、Deep ThinkはARC-AGI-2において84.6%という記録的なスコアを叩き出し(Gemini 3.1 Proの77.1%をさらに上回る)、事実上、特定領域において人間レベルの抽象推論に肉薄している<sup>1</sup>。さらに、競技プログラミングのCodeforcesでは3,455 Eloという驚異的な数値を記録し、学术界の最難関テストとされる「Humanity's Last Exam(人類最後の試験)」においても、ツールなしで44.4%、ツールありで51.4%という前人未達のスコアを達成した<sup>1</sup>。また、物理学オリンピック(87.7%)や化学オリンピック(82.8%)においても金メダルレベルの定量的推論

能力を示している<sup>1</sup>。

## 5.2. トレードオフの克服:「3段階の思考深度システム(Deep Think Mini)」

しかし、この並列探索と自己検証のループは、APIのレスポンス時間(レイテンシ)の大幅な遅延と、クエリあたりの計算コスト(コンピュート予算)の劇的な増加を引き起こすという避けられないトレードオフをもたらす<sup>1</sup>。数学の証明や新規アルゴリズムの発見、数日にわたる自律型リサーチシステム(Gemini Deep Research 2.0など)には最適だが、日常的なコーディング支援、リアルタイムのシステム統合、高速なデータパイプライン処理には過剰なオーバースペックとなる<sup>1</sup>。

ここで重要になるのが、Gemini 3.1 Proの立ち位置である。Gemini 3.1 Proは、Deep Thinkの高度な推論パイプラインで培われたコア・インテリジェンスの「エッセンス」を抽出し、より実用的な速度とコストに最適化して組み込んだモデルとして設計されている<sup>2</sup>。

この実用化を支える最大の機能が、Gemini 3.1 ProのAPIにおいて新たに導入された\*\*「3段階の思考深度システム(Deep Think Mini)」\*\*である<sup>11</sup>。開発者はタスクの性質に応じて、APIのパラメータ(thinking\_level)を通じて、以下の3つのレベルでエージェントの「推論予算(Reasoning Budget)」を柔軟に制御できる<sup>11</sup>。

思考深度レベル	特徴と推論能力	最適なユースケース	レイテンシへの影響
<b>HIGH</b>	現在のDeep Thinkの「ミニ版」に近い深い推論。マルチステップ計画と自己検証を伴う。	複雑なデバッグ、数学的証明、エージェントの自律的な戦略立案、検証済みコード生成。	高い(レスポンスに時間を要する)
<b>MEDIUM</b>	前世代(3.0 Pro)の最高レベル(high)に相当する推論品質。バランスの取れたアプローチ。	コードレビュー、技術分析、アーキテクチャ設計、日常的なプログラミング支援。	中程度(適度な速度と精度)
<b>LOW (MINIMAL)</b>	推論コストを最小限に抑え、事前の計画プロセスを省略して高速な生成を優先する。	データ抽出、フォーマット変換、単純な質疑応答、翻訳タスク。	極めて低い(瞬時に応答)

注: Gemini 3.1 Proでは推論プロセス自体を完全に無効化(オフ)にすることはできず、最低でもLOW

(あるいはMINIMAL)レベルの思考シグネチャが要求される<sup>29</sup>。

開発者は、過去において「Gemini 3.0 Proのhighモード」で実行していたタスクを、3.1 Proでは「MEDIUMモード」に切り替えるだけで、同等以上の推論品質をより高速に得ることができる<sup>8</sup>。真に深い推論が必要な局面でのみ「HIGHモード」を活用することで、コストとパフォーマンスの最適化をアプリケーションレベルで実現可能となった。

### 5.3. エージェント開発プラットフォーム「Google Antigravity」との統合

Gemini 3.1 Proのこれらの推論能力とエージェント機能は、Googleが新たに提供を開始したエージェントファーストの開発プラットフォーム「Google Antigravity」を通じて最大限に発揮される<sup>5</sup>。

Antigravityは、単なるコードエディタの延長(オートコンプリートツール)ではなく、開発者が複数の自律型エージェントをスポン(生成)し、オーケストレーションするための専用環境である<sup>32</sup>。Gemini 3.1 Proを搭載したAntigravity上のエージェントは、エディタ、ターミナル、そしてブラウザを横断して自律的に計画、実行、検証を行う<sup>32</sup>。

例えば、レガシーな認証モジュールのリファクタリングや、巨大なデータベース移行といった長時間のタスクをエージェントに委任することができる<sup>31</sup>。エージェントは自律的にコードを書き換え、ターミナルを用いてテスト環境を立ち上げ、ブラウザを通じてUIの変更を確認する<sup>32</sup>。このプロセスにおいて、開発者は膨大な生ログを追う必要はない。Gemini 3.1 Proは「アーティファクト(タスクリスト、実装計画書、スクリーンショット、ブラウザの操作録画など)」と呼ばれる有形の成果物を生成し、開発者はそれらを通じてエージェントの論理プロセスを一目で検証することができる<sup>32</sup>。修正が必要な場合は、ドキュメントにコメントを残す感覚でアーティファクトに直接フィードバックを与えれば、エージェントは実行フローを停止することなく指示を統合し、作業を継続する<sup>32</sup>。

前述のcustomtoolsエンドポイントやthinking\_levelの柔軟な制御は、このAntigravityプラットフォームの裏側でシームレスに機能しており、長時間の自律タスクにおけるモデルの脱線(幻覚やループ)を極限まで抑え込んでいる<sup>8</sup>。

## 6. 結論と戦略的提言

Gemini 3.1 Proは、単なるAIモデルのマイナーチェンジの枠を大きく超え、汎用人工知能(AGI)に向けた「推論時間計算量のスケーリング」がいかに実用的なフェーズに移行したかを示す象徴的なプロダクトである。ARC-AGI-2における77.1%という圧倒的なスコアは、AIが単なる「知識の検索・要約エンジン」から、未知のパターンを解釈し独自の論理法則を導き出す「適応型・流動的推論エンジン」へと変貌を遂げたことを歴史的に証明している。

Claude Opus 4.6やGPT-5.3 Codexといった他社のフロンティアモデルと比較した場合、極度に深いドメイン特化型のバグ修正や、ターミナルでの低レイヤ自律操作では一部一歩譲る領域が存在することは事実である。また、100万トークンの超長文脈入力時の情報抽出精度の低下というアーキテクチャ上の弱点も残されている。

しかしながら、LiveCodeBench Proで実証された最高峰のアルゴリズム構築能力、複雑な複数ツ

ルのオーケストレーション能力、革新的な「3段階の思考深度システム (Deep Think Mini)」による推論予算の緻密なコントロール、そして競合他社を圧倒する経済性(7.5倍のコスト優位性)の組み合わせにおいて、Gemini 3.1 Proは現在最もバランスが取れ、企業規模での実運用に向けた最適解であると結論付けられる。

組織のAI戦略においては、Gemini 3.1 Proと特殊推論モードであるGemini Deep Thinkを適材適所で使い分ける「インテリジェント・モデル・ルーティング」の構築が急務となる。日常的なAPI処理やUI生成はGemini 3.1 ProのLOW/MEDIUMモードで高速処理し、複雑なエージェントワークフローにはHIGHモードやcustomtoolsエンドポイントを適用、そして極限の精度が求められる研究開発や金融モデリングにおいてのみDeep Thinkモードへエスカレーションする階層型アーキテクチャが求められる。

今後、Google Antigravityのようなプラットフォームの普及に伴い、Gemini 3.1 Proは「コードを書くAI」から「複雑な業務システムを自律的に設計・統合・検証する知能のオーケストレーター」へと、ソフトウェアエンジニアリングと企業プロセスのあり方を根本から再定義していくであろう。

## 引用文献

1. Gemini 3 Deep Think: Reasoning Benchmarks & Complete Guide, 2月 20, 2026にアクセス、  
<https://www.digitalapplied.com/blog/gemini-3-deep-think-reasoning-benchmarks-guide>
2. Google releases Gemini 3.1 Pro: Here's what's new and who gets it first, 2月 20, 2026にアクセス、  
<https://timesofindia.indiatimes.com/technology/tech-news/google-releases-gemini-3-1-pro-heres-whats-new-and-who-gets-it-first/articleshow/128569493.cms>
3. Gemini 3.1 Pro Complete Guide 2026: Benchmarks, Pricing, API, 2月 20, 2026にアクセス、  
<https://www.nxcode.io/en/resources/news/gemini-3-1-pro-complete-guide-benchmarks-pricing-api-2026>
4. Techmeme, 2月 20, 2026にアクセス、<https://www.techmeme.com/260219/p30>
5. Gemini 3.1 Pro の概要 | npaka - note, 2月 20, 2026にアクセス、  
<https://note.com/npaka/n/nc859d1b06094>
6. Gemini 3.1 Pro: A smarter model for your most complex tasks, 2月 20, 2026にアクセス、  
<https://blog.google/innovation-and-ai/models-and-research/gemini-models/gemini-3-1-pro/>
7. 【徹底解説】Gemini 3.1 Pro 登場。実力を4モデル比較で検証 | AGI ..., 2月 20, 2026にアクセス、  
<https://chatgpt-lab.com/n/n64d9f440c178>
8. Gemini 3.1 Pro vs 3.0 Pro Preview Full Comparison, 2月 20, 2026にアクセス、  
<https://help.apiyi.com/en/gemini-3-1-pro-vs-3-pro-preview-comparison-guide-en.html>
9. Gemini 3.1 Pro - Model Card - Google DeepMind, 2月 20, 2026にアクセス、  
<https://deepmind.google/models/model-cards/gemini-3-1-pro/>
10. Gemini 3 Developer Guide | Gemini API - Google AI for Developers, 2月 20, 2026

- にアクセス、<https://ai.google.dev/gemini-api/docs/gemini-3>
11. Gemini 3.1 Pro Preview APIがAPIYIでリリース: 推論性能を2倍に ..., 2月 20, 2026にアクセス、  
<https://help.apiyi.com/ja/gemini-3-1-pro-preview-api-available-apiyi-guide-ja.html>
  12. Gemini 3.1 Pro on Gemini CLI, Gemini Enterprise, and Vertex AI, 2月 20, 2026にアクセス、  
<https://cloud.google.com/blog/products/ai-machine-learning/gemini-3-1-pro-on-gemini-cli-gemini-enterprise-and-vertex-ai>
  13. Gemini 3.1 Pro - Google DeepMind, 2月 20, 2026にアクセス、  
<https://deepmind.google/models/gemini/pro/>
  14. AI Gemini 3.1 Pro Hadir, Google Tantang Dominasi GPT dan Claude, 2月 20, 2026にアクセス、  
<http://tekno.kompas.com/read/2026/02/20/09450067/ai-gemini-3.1-pro-hadir-google-tantang-dominasi-gpt-dan-claude>
  15. My 2026 AI Predictions: Agents Get Real, Benchmarks Get Weird, 2月 20, 2026にアクセス、  
<https://adam.holter.com/my-2026-ai-predictions-agents-get-real-benchmarks-get-weird-and-continual-learning-stays-external/>
  16. Google's new AI model with double the reasoning power - Xpert.Digital, 2月 20, 2026にアクセス、  
<https://xpert.digital/en/google-gemini-3.1-pro/>
  17. ARC-AGI-2: A New Challenge for Frontier AI Reasoning Systems, 2月 20, 2026にアクセス、  
<https://arxiv.org/abs/2505.11831>
  18. GPT-5.2 & ARC-AGI-2: A Benchmark Analysis of AI Reasoning, 2月 20, 2026にアクセス、  
<https://intuitionlabs.ai/articles/gpt-5-2-arc-agi-2-benchmark>
  19. ARC-AGI-2, 2月 20, 2026にアクセス、  
<https://arcprize.org/arc-agi/2/>
  20. Gemini 3 Flash: Pricing, Context Window, Benchmarks, and More, 2月 20, 2026にアクセス、  
<https://llm-stats.com/models/gemini-3-flash-preview>
  21. ARC Prize launches its toughest AI benchmark yet: ARC-AGI-2, 2月 20, 2026にアクセス、  
<https://www.artificialintelligence-news.com/news/arc-prize-launches-toughest-ai-benchmark-yet-arc-agi-2/>
  22. Gemini 3 Deep Think Achieves 45.1% on ARC-AGI-2 - Reddit, 2月 20, 2026にアクセス、  
[https://www.reddit.com/r/accelerate/comments/1p0go5r/gemini\\_3\\_deep\\_think\\_achieves\\_451\\_on\\_arcagi2/](https://www.reddit.com/r/accelerate/comments/1p0go5r/gemini_3_deep_think_achieves_451_on_arcagi2/)
  23. Google Gemini 3 Benchmarks (Explained) - Vellum, 2月 20, 2026にアクセス、  
<https://www.vellum.ai/blog/google-gemini-3-benchmarks>
  24. Gemini 3の思考モード・Proモード・高速モードの違い、上限を解説！, 2月 20, 2026にアクセス、  
<https://www.ai-souken.com/article/gemini-3-modes-comparison>
  25. A new era of intelligence with Gemini 3 - Google Blog, 2月 20, 2026にアクセス、  
<https://blog.google/products-and-platforms/products/gemini/gemini-3/>
  26. Is This AGI? Google's Gemini 3 Deep Think Shatters Humanity's Last, 2月 20, 2026にアクセス、  
<https://www.marktechpost.com/2026/02/12/is-this-agi-googles-gemini-3-deep-t>

[hink-shatters-humanitys-last-exam-and-hits-84-6-on-arc-agi-2-performance-to-day/](#)

27. Gemini 3 vs Gemini Pro vs Gemini DeepThink: Key Differences, 2月 20, 2026にアクセス、<https://tech-now.io/en/blogs/gemini-3-vs-gemini-pro-vs-gemini-deepthink>
28. Gemini 3 Deep Think vs. Gemini Deep Research - SourceForge, 2月 20, 2026にアクセス、  
<https://sourceforge.net/software/compare/Gemini-3-Deep-Think-vs-Gemini-Deep-Research/>
29. Thinking | Generative AI on Vertex AI - Google Cloud Documentation, 2月 20, 2026にアクセス、  
<https://docs.cloud.google.com/vertex-ai/generative-ai/docs/thinking>
30. Gemini thinking | Gemini API - Google AI for Developers, 2月 20, 2026にアクセス、  
<https://ai.google.dev/gemini-api/docs/thinking>
31. gemini-3-1-pro-in-google-antigravity, 2月 20, 2026にアクセス、  
<https://antigravity.google/blog/gemini-3-1-pro-in-google-antigravity>
32. Build with Google Antigravity, our new agentic development platform, 2月 20, 2026にアクセス、  
<https://developers.googleblog.com/build-with-google-antigravity-our-new-agentic-development-platform/>