

Gemini 3.5 Flashのエージェント性能向上とIP業務への示唆

作成者: Manus AI

作成日: 2026年5月23日

主題: 「Gemini 3.5 FlashのGDPval-AA Elo 1656が、複数エージェントによる特許調査・先行技術検索ワークフローの実用可能性を高める」という見方の詳細調査

1. 要旨

Gemini 3.5 Flashは、Google DeepMindが2026年5月19日に公表したGemini 3系列の高速・マルチモーダル推論モデルであり、公式モデルカードでは最大**1Mトークンの入力テキスト**、**64Kトークンのテキスト出力**、テキスト・画像・音声・動画入力への対応が示されている。^① Googleは同モデルを、**agentic workflows**、**coding tasks**、**multi-week enterprise processes**に適したモデルとして位置付けている。^①

特に注目されるのは、Gemini 3.5 Flashが**GDPval-AAでElo 1656**を記録した点である。GDPval-AAはArtificial AnalysisがOpenAIのGDPvalデータセットを用いて実施するエージェント評価であり、モデルにWebブラウジングとシェルアクセスを与え、44職種・9主要産業にまたがる実世界型タスクを解かせ、ブラインドのペアワイズ比較からEloを算出する。^② この指標は特許調査そのものを測るものではないが、文書・図表・スプレッドシート等を含む実務成果物の生成能力、ツール利用能力、反復的な作業遂行能力を反映するため、**先行技術検索**、**クレーム要素分解**、**文献スクリーニング**、**クレームチャート作成**、**調査報告書作成**といったIP業務との関連性が高い。

ただし、結論は慎重に読む必要がある。GDPval-AA 1656は、Gemini 3.1 Proの1314やGemini 3 Flashの1204から大きく改善している一方で、同リーダーボード上ではGPT-5.5やClaude Opus 4.7等の一部高努力設定を下回る。^① ^② したがって、Gemini 3.5 Flashの意義は「全タスクで最高精度」というよりも、**高いエージェント性能を高速・低遅延・比較的lowコストで実行できる可能性**にある。IP実務では、これを単独モデルとしてではなく、検証済み特許データベース、検索API、引用根拠、ログ、専門家レビューを組み合わせた**監督付きマルチエージェント・ワークフロー**として設計することが重要である。

2. Gemini 3.5 Flashの性能指標とIP業務への関係

Google DeepMindのモデルカードによれば、Gemini 3.5 Flashはエージェント、ツール利用、UI操作、専門タスク、マルチモーダル理解、長文脈処理に関する複数の評価で高いスコアを示している。^① IP業務に直接関係するのは、GDPval-AAだけでなく、MCP Atlas、Toolathlon、OSWorld、Finance Agent v2、CharXiv Reasoning、MMMU-Pro、長文脈評価である。特許調

査は単なる自然言語応答ではなく、外部検索、文献読解、表・図面の確認、引用根拠の抽出、成果物化を含むため、これらの複合的な指標が重要になる。

評価軸	Gemini 3.5 Flashの公表値	IP業務との関連	解釈
GDPval-AA	Elo 1656	知識労働型タスク、実務成果物生成	特許調査報告書、分析メモ、クレームチャート等に近い作業形式を含む可能性があるが、特許専用評価ではない。
MCP Atlas	83.6%	ツール連携・マルチステップ処理	特許DB、検索API、社内DMS、表計算、ナレッジベースを連携するワークフローに関係する。
Toolathlon	56.5%	一般的な現実世界ツール利用	検索、ファイル処理、表作成、外部サービス連携を伴うIP調査に関係する。
OSWorld-Verified	78.4%	UI操作・コンピュータ利用	特許検索プラットフォームやWeb UIの操作自動化に関係する。
CharXiv Reasoning	84.2%	図表・チャートからの情報合成	特許図面、フローチャート、実験データ表、化学・機械図面の読解に関係する。
入力コンテキスト	最大1Mトークン	長大な特許明細書・審査履歴・引用文献束	複数文献を同時に扱う調査や、審査経過分析に有利。ただし1M評価では性能低下も報告されている。

Google公式ブログは、3.5 Flashを「複雑なエージェント型ワークフローの実行支援」に向けたモデルと説明し、Antigravity harnessと組み合わせることで、複数サブエージェントを監督下で動かす長期・複雑タスクに適すると述べている。³ IP業務の観点では、この説明は、単一のチャットボットではなく、**専門機能を持つ複数エージェントの協調**によって調査品質を高める方向性と整合する。

Googleは、Gemini 3.5 Flashについて、複雑な長期タスクで実用性を発揮し、Antigravity

harnessと組み合わせることで協調的なサブエージェントを展開できると説明している。③

3. GDPval-AA Elo 1656の意味

GDPval-AAは、OpenAIが提案したGDPvalをArtificial Analysisがエージェント評価として実装したものである。OpenAIのGDPvalは、米国GDPへの産業寄与を基に選定した9産業・44職種を対象とし、フルセット1,320タスク、公開gold set 220タスクから構成される。タスクは平均14年以上の経験を持つ専門家が作成・レビューし、法的メモ、設計図、顧客対応、看護計画など、実際の業務成果物に近い形式を含む。④

Artificial AnalysisのGDPval-AAは、このGDPvalデータセットに対して、モデルにシェルアクセスとWebブラウジングを与えるエージェント・ハーネス「Stirrup」を用いる。Eloは、モデル出力同士のブラインド・ペアワイズ比較から導出される。② そのため、GDPval-AAは従来型の閉じたQ&Aベンチマークよりも、**検索、情報統合、成果物作成、ツール利用**に近い評価だと考えられる。

観点	GDPval	GDPval-AA	IP業務への読み替え
対象	44職種・9産業の実務タスク	GDPvalをエージェントハーネスで実行	特許調査は「法律・技術・文書分析」が融合した知識労働であり、近縁性がある。
成果物	文書、スライド、図表、表計算等	ツール利用を含む成果物生成	先行技術調査報告、無効資料調査、FTO分析、クレームチャートに近い。
評価方法	専門家/比較評価	ブラインドペアワイズ比較からElo	相対評価であり、個別案件での法的妥当性を直接保証しない。
制限	初期版はone-shot評価の制約あり	特許専用ではない	IP実務では反復検索、専門家レビュー、管轄法・審査基準への適合検証が必須。

重要なのは、GDPval-AA 1656が**特許調査の完全自動化を意味しない**ことである。OpenAI自身もGDPvalについて、初期版はone-shot評価であり、実務の反復性、曖昧性、複数ドラフトでの改善を十分に反映しないと説明している。④ 特許調査では、検索式の改善、ノイズ除去、追加観点の発見、出願人・発明者・CPC・引用関係・ファミリー調査等を反復するため、この制約は大きい。

4. 特許調査・先行技術検索がエージェント化に向く理由

先行技術検索は、情報検索、専門的読解、法的判断補助、成果物作成が結合したワークフローである。WIPO Magazineは、特許庁において出願件数と技術複雑性が増大し、JPOの2018年調査では審査官が時間の30%を先行技術検索、10%をその理解に費やすと紹介している。⁵ また、WIPOは各国IP庁のAI活用状況を整理しており、CIPOが商用の意味検索エンジンを先行技術・引用検索に用いていること、IP Australiaが自動予備検索や外国審査レポート分析を進めていることを示している。⁶

特許検索レビュー研究も、先行技術検索の困難性として、文献量の増大、専門的・抽象的な特許語彙、言語・法域差、検索クエリ設計の難しさ、重要文献の見落としリスクを挙げている。⁷ この構造は、エージェント型AIが得意とされる「タスク分解」「反復探索」「外部ツール利用」「結果の統合」に適している。

特許調査の工程	従来の主な課題	エージェント化で期待される支援
発明把握	技術的特徴と法的なクレーム要素の切り分けが難しい	発明開示・請求項を構成要素に分解し、検索観点を複数生成する。
検索戦略設計	キーワード、同義語、分類、引用、ファミリの組合せに熟練が必要	主要特徴ごとに検索式、CPC/IPC、類義語、英日中表現を提案する。
先行技術検索	法域・言語・非特許文献の分散、検索漏れ	特許DB・学術DB・Web検索を役割分担して反復検索する。
スクリーニング	大量ノイズと候補文献の優先順位付け	要素一致度、日付、法域、引用関係で候補をランキングする。
クレーム対応付け	どの文献のどの段落がどの要素を開示するか確認が重い	引用箇所付きでクレームチャート案を作成する。
品質保証	ハルシネーション、引用誤り、過大評価	検証エージェントが特許番号、公開日、引用段落、根拠の有無を監査する。

実際、PatExpertという研究では、メタエージェントが特許分類、受理予測、請求項生成、要約、複数特許分析、科学的仮説生成などの専門エージェントを統括し、批評エージェントが出力を評価・改善するマルチエージェント構成が提案されている。⁸ これは、Gemini 3.5 Flashのような高速エージェントモデルを用いる場合の実装パターンとして参考になる。

5. 実用ワークフロー案: 複数エージェントによる先行技術検索

Gemini 3.5 FlashをIP業務に活用する場合、推奨される設計は、単一モデルに「先行技術を探して」と依頼する形ではない。PQAIの実験的解説が示すように、ChatGPT等の汎用LLM単体で特許番号や文献を生成させると、実在するが無関係な特許番号を提示するなど、高リスクのハルシネーションが起り得る。⁹ したがって、LLMは必ず、検証済み特許データベース、意味検索API、引用取得機構、ログ保存、専門家レビューと組み合わせる必要がある。

実務導入時には、次のような**監督付きマルチエージェント構成**が現実的である。

エージェント	役割	主な入力	主な出力	必須ガードレール
発明解析エージェント	発明の課題、構成、効果、必須要素を抽出	発明提案書、請求項案、図面	要素分解表、検索観点	人間が必須要素を承認する。
検索戦略エージェント	キーワード、分類、同義語、検索式を生成	要素分解表	検索式セット、CPC/IPC候補	検索式と除外条件を記録する。
データベース検索エージェント	特許DB・非特許文献DBを検索	検索式、分類	候補文献リスト	LLM生成文献ではなくDB取得文献のみ採用する。
スクリーニングエージェント	候補文献を要素一致で順位付け	候補文献全文/要約	優先文献、除外理由	根拠箇所がない一致判定を禁止する。
クレームチャートエージェント	請求項要素と文献箇所を対応付け	請求項、候補文献	要素別対応表	段落番号・請求項番号・図番を必須化する。
反証/品質監査エージェント	見落とし、引用誤り、過大評価を検出	全ログ、検索結果、チャート	監査コメント、追加検索案	別検索戦略で再検索する。
レポート作成エージェント	調査報告書を作成	チャート、監査結果	報告書案	最終判断は弁理士・専門家が行う。

この構成では、Gemini 3.5 Flashの高速性は、検索式の多案生成、候補文献の初期スクリーニング、要素対応表の作成、監査サイクルの短縮に効く可能性がある。特に、複数エージェントを並行的または反復的に動かす場合、速度とコストは品質と同じくらい重要な制約となる。

Googleおよび解説記事は、3.5 Flashが他のフロンティアモデルより高速で、長期エージェントワークフローのコスト面で有利になり得ると説明している。^{3 10}

6. 具体的なユースケース別評価

ユースケース	実用可能性	Gemini 3.5 Flashが寄与し得る点	注意点
発明提案段階の予備的な新規性調査	高い	自然文の発明説明をクレーム要素に分解し、検索観点を複数提示できる。	網羅調査ではなく、初期スクリーニングとして位置付けるべきである。
出願前先行技術調査	中～高	長文脈とツール利用により、複数文献の比較とレポート化を支援できる。	最終的な新規性・進歩性判断は専門家レビューが必要である。
無効資料調査	中	請求項要素に対する文献対応付けと検索観点拡張に有用。	見落としリスクが高く、複数DB・専門検索者の併用が必須である。
FTO調査	中	製品特徴と有効特許請求項の対応付けを補助できる。	権利範囲解釈、均等論、法域別リスク評価は人間判断が不可欠である。
OA対応・拒絶理由分析	高い	引用文献と請求項補正案の比較、反論骨子作成を支援できる。	代理人の法的判断、包袋履歴、禁反言リスクの確認が必要である。
特許ランドスケープ	高い	大量文献のクラスタリング、要約、技術トレンド抽出に適する。	定量分析ではデータ正規化と母集団定義が結果を左右する。

7. リスクとガバナンス

IP業務では、AI活用による効率化の一方で、ハルシネーション、秘匿情報漏えい、根拠不明な法的判断、検索漏れ、データベース利用規約違反、発明者・代理人・企業秘密の取り扱いが重大なリスクとなる。特に、特許調査では「それらしく見えるが無関係な文献」や「存在するが技術内容が異なる特許番号」は実害が大きい。したがって、Gemini 3.5 Flashのような高性能モデルであっても、**根拠文献の实在確認、引用箇所確認、検索ログ保存、専門家レビュー、機密データ管理**を必須要件にするべきである。

リスク	発生例	推奨対策
-----	-----	------

ハルシネーション	架空または無関係な特許番号を提示する	DB/APIで取得した文献のみを候補にし、出典URL・公報番号・段落番号を必須化する。
検索漏れ	同義語・別分類・外国語文献を見落とす	複数検索戦略、CPC/IPC、引用ネットワーク、ファミリー検索、非特許文献検索を組み合わせる。
過大な法的判断	AIが新規性・進歩性を断定する	AI出力は「調査補助・論点整理」とし、最終判断は弁理士・弁護士・知財専門家が行う。
秘密情報漏えい	未公開発明を外部モデル/APIに入力する	エンタープライズ契約、データ保持設定、匿名化、オンプレ/閉域検索環境を確認する。
再現性不足	どの検索式・候補を用いたか追跡できない	プロンプト、検索式、検索日時、DB、除外理由、モデル設定をログ化する。

8. 導入ロードマップ

実務導入は、いきなり全面自動化を狙うよりも、既存の特許調査プロセスに段階的に組み込むのが望ましい。最初の対象は、リスクが比較的低く、成果検証が容易な**予備調査、技術分類、候補文献要約、OA引用文献整理**が適している。その後、既知の審査引用文献や過去の無効資料調査案件をベンチマークとして、再現率、精度、時間短縮、専門家修正量を評価するべきである。

フェーズ	期間目安	実施内容	成功指標
PoC	2~4週間	過去案件10~20件で、検索観点生成・要約・候補順位付けを試す。	既知重要文献の再発見率、専門家修正時間、ハルシネーション件数。
パイロット	1~3か月	特定技術分野でマルチエージェント調査フローを運用する。	検索時間削減、候補文献の精度、レビュー負荷、利用者満足度。
本番化	3~6か月	特許DB/API、社内DMS、ログ、権限管理、監査機能と統合する。	再現性、監査可能性、セキュリティ、標準業務への定着。
高度化	継続	クレームチャート、OA対応、FTO、ランドス	品質指標、コスト削減、リスク低減、専門

9. 結論

Gemini 3.5 FlashのGDPval-AA Elo 1656は、IP業務において、特に**複数エージェントによる特許調査・先行技術検索ワークフロー**の実用可能性を高める重要なシグナルである。GDPval-AAは、Web・シェル利用を含むエージェント型の実世界知識労働を評価しており、単なる会話性能よりも、特許調査のような「検索・分析・根拠提示・成果物化」を伴う業務に近い。

一方で、同スコアは特許専用の正確性、再現率、法的妥当性を直接保証するものではない。したがって、実務上の最適解は、Gemini 3.5 Flashを単独の判断主体として使うのではなく、**検証済みデータベース、意味検索API、クレーム要素分解、根拠抽出、反証エージェント、専門家レビューを組み合わせた監督付きAI調査基盤**を構築することである。これにより、先行技術検索の探索範囲を広げ、初期スクリーニングとクレームチャート作成を高速化し、知財担当者・弁理士がより高度な判断に集中できる可能性がある。

References

- [1] Google DeepMind, Gemini 3.5 Flash Model Card
- [2] Artificial Analysis, GDPval-AA Leaderboard
- [3] Google Blog, Gemini 3.5: frontier intelligence with action
- [4] OpenAI, Measuring the performance of our models on real-world tasks
- [5] WIPO Magazine, Patent Office Sustainability and the Role of Artificial Intelligence
- [6] WIPO, Index of AI initiatives in IP offices
- [7] Ali et al., Innovating Patent Retrieval: A Comprehensive Review of Techniques, Trends, and Challenges in Prior Art Searches
- [8] Srinivas et al., Towards Automated Patent Workflows: AI-Orchestrated Multi-Agent Framework for Intellectual Property Management and Analysis
- [9] PQAI, How to Get Better Prior Art Results in ChatGPT With PQAI Integration?
- [10] DataCamp, Gemini 3.5 Flash: Google's Fastest Agentic Model
- [11] USPTO, USPTO launches new AI Pilot for pre-examination utility application search
- [12] USPTO, Artificial Intelligence Search Automated Pilot Program
- [13] Ikoma and Mitamura, Can AI Examine Novelty of Patents?
- [14] Risch et al., PatentMatch: A Dataset for Matching Patent Claims & Prior Art
- [15] Lo et al., Large Language Model Informed Patent Image Retrieval
- [16] VentureBeat, Google says Gemini 3.5 Flash can slash enterprise AI costs...