

# Gemma 4 12Bが知財業務にもたらす変革と 将来展望：完全ローカル・マルチモーダルAIに よる実務プロセスの再構築

Gemini 3.1 pro

## 1. イントロダクション：知財実務における生成AIの進化とセキュリティのジレンマ

知的財産 (IP) 業務、とりわけ特許出願、先行技術調査、および拒絶理由通知への対応 (Office Action 応答) といった実務は、極めて高度な論理的思考力、広範な技術理解、そして厳密な機密保持が要求される専門領域である。近年、大規模言語モデル (LLM) を中心とする生成AIの進化は、これらの高度なナレッジワークの生産性を飛躍的に向上させる可能性を示してきた。しかしながら、知財業界は常に「利便性とセキュリティのトレードオフ」という根本的なジレンマに直面してきた。特許出願前の発明に関する未公開情報は、企業にとって最高の機密情報 (営業秘密) である。これをクラウドベースのパブリックなSaaS型生成AI (例えば一般的なAPI経由での利用) に入力することは、外部サーバーへのデータ送信を通じた情報漏洩や、AIの学習データとして利用されることによる特許法上の「新規性喪失」という致命的なリスクを伴うためである<sup>1</sup>。

このような背景の中、2026年6月3日、Googleは新たなオープンウェイトモデル「Gemma 4 12B」をリリースした<sup>3</sup>。このモデルは、119.5億 (12B) のパラメータを持ちながら、データセンター級の専用AIハードウェアを必要とせず、一般的なエンタープライズ向けのコンシューマーノートPC (16GBのVRAMまたはユニファイドメモリ) 上で完全にローカル稼働するように設計されている点が最大の長である<sup>3</sup>。単なる軽量化モデルの枠を超え、テキスト、画像、音声ネイティブに処理するエンコーダフリーの「統合アーキテクチャ (Unified architecture)」を採用し、256Kトークンという長大なコンテキストウィンドウを備えている<sup>3</sup>。

本レポートでは、Gemma 4 12Bの技術的特性とアーキテクチャを詳細に分析し、それが特許事務所や企業の知財部門における業務プロセス (発明発掘から明細書作成、権利化対応まで) にどのようなパラダイムシフトをもたらすのかを、技術的、実務的、およびコンプライアンスの観点から網羅的に論じる。

## 2. Gemma 4 12Bの概要とアーキテクチャの革新性

Gemma 4 12Bの影響を正確に評価するためには、まずその技術的基盤と、旧世代のモデルからどのように進化したのかを理解する必要がある。

### 2.1. Gemma 4ファミリーと12Bモデルの位置づけ

Gemma 4ファミリーは、高度な推論とエージェント的ワークフローの構築を目的として設計されたオープンモデル群であり、Apache 2.0ライセンスの下で提供されている<sup>3</sup>。このファミリーは、用途とハードウェアの制約に応じて複数のサイズで構成されており、それぞれがエッジデバイスからハイエンドのワークステーション、サーバーまで幅広い環境にデプロイできるようスケラブルな設計となっ

ている<sup>9</sup>。

以下の表は、Gemma 4ファミリーを構成する主要モデルの仕様比較である<sup>8</sup>。

プロパティ	E2B	E4B	12B	26B A4B MoE	31B Dense
合計パラメータ数	2.3B (有効: エンベディング含み5.1B)	4.5B (有効: エンベディング含み8B)	約12B	25.2B (総) / 3.8B (有効)	30.7B
レイヤー数	35	42	未公開	30	60
コンテキスト長	128Kトークン	128Kトークン	256Kトークン	256Kトークン	256Kトークン
サポートされるモダリティ	テキスト, 画像, 音声	テキスト, 画像, 音声	テキスト, 画像, 音声	テキスト, 画像	テキスト, 画像
推論機構 / Expert数	Dense	Dense	Dense	8 active / 128 total (1 shared)	Dense

このファミリーの中で、12Bモデルは、エッジデバイス向けのE4B(Effective 4Bパラメータ)と、より大規模な26B MoE (Mixture-of-Experts) の間に位置する「スイートスポット」として機能する<sup>13</sup>。26B MoEモデルは、総パラメータ数25.2Bでありながら推論時には3.8Bのパラメータのみをアクティブにする(8/128 Expert)ことで高速化を図っているが<sup>8</sup>、それでも一定以上のVRAMを要求する。これに対し、12Bモデルは、量子化認識トレーニング(QAT: Quantization-Aware Training)などの最適化によりメモリ要件を劇的に削減しており、16GBのRAMを搭載した標準的なノートPCでスタンドアロン稼働する<sup>3</sup>。これにより、弁理士や知財担当者が日常的に使用するモバイルワークステーション上で、クラウドに匹敵する「フロンティアインテリジェンス」を完全にオフラインで実現できる環境が整った<sup>10</sup>。

## 2.2. エンコーダフリー「統合アーキテクチャ」がもたらすマルチモーダル処理の飛躍

従来のマルチモーダルAIは、テキストを処理する主要な言語モデル(LLMバックボーン)に加えて、音声波形や画像データをLLMが理解できる表現に変換するための専用の「エンコーダ(視覚エンコーダや音声エンコーダ)」を個別に統合する構造を持っていた。この分離された構造は、推論プロセスにおける遅延(レイテンシ)を増加させ、システム全体のメモリ使用量を肥大化させるという構造的な課題を抱えていた<sup>3</sup>。

Gemma 4 12Bにおける最も根本的な技術的ブレイクスルーは、この視覚および音声エンコーダを完

全に排除した「エンコーダフリー(Encoder-Free)」な統合アーキテクチャの採用である<sup>3</sup>。Google DeepMindの開発チームは、視覚情報の処理を、単一の行列乗算、位置埋め込み、および正規化からなる軽量な埋め込みモジュールに置き換えた。これにより、視覚的なパッチ情報はLLMのバックボーンに直接送られ、言語モデル自体が視覚処理を主導する仕組みとなっている<sup>3</sup>。さらに音声処理に関しては、音声エンコーダを完全に削除し、生の音声波形(オーディオ信号)を直接テキストトークンと同じ次元空間に投影するという画期的なアプローチを採用している<sup>3</sup>。

知財実務において、このアーキテクチャの変更は極めて実践的な価値をもたらす。特許明細書は、高度な技術文章(テキスト)と、フローチャート、回路図、機械の構造図などの「図面」が不可分に結びついたハイブリッドな情報媒体である。Gemma 4 12Bは、OCR(光学文字認識)やチャート理解において可変アスペクト比と可変解像度をネイティブにサポートしている<sup>8</sup>。さらに、画像の視覚的詳細度(解像度)をコントロールするためのビジュアルトークン予算(70, 140, 280, 560, 1120トークン等の段階)を構成可能であり、粗い図面は低演算で高速に、微細な半導体回路図などは高解像度で精緻に処理するという柔軟な運用が可能である<sup>11</sup>。これにより、明細書本文中の「モーター(10)の回転軸(12)が...」というテキスト記述と、添付された図面上の該当符号を、ローカル環境の限られた計算資源の下でもシームレスかつ高速に紐付けて理解することが可能となる。

### 3. 知財業界におけるセキュリティ課題とローカルAIの台頭

Gemma 4 12Bが知財業界において「ゲームチェンジャー」と目される理由は、純粋な性能の高さだけでなく、業界特有の極めて厳格なコンプライアンスとセキュリティ要件に完璧に合致するアーキテクチャを備えているからである。

#### 3.1. クラウドAIの新規性喪失リスクとオンプレミスの限界

特許法制において、出願前に発明の内容が公知となった場合、原則として新規性を喪失し、特許を受けることができなくなる(新規性喪失の例外規定等の救済措置は存在するが、極めて限定的である)。そのため、開発段階の未出願の発明情報を、APIを通じて外部のクラウドサービスに送信することに対しては、知財部門から強い懸念が示されてきた<sup>1</sup>。特に、入力したデータがLLMプロバイダーの将来のモデル学習(再学習)に利用される可能性は、他社への情報漏洩と同義であり、企業法務において絶対に回避すべきリスクとされている。

この課題に対処するため、業界では二つのアプローチが取られてきた。第一は、学習にデータを利用しないエンタープライズ向けクラウドサービスの利用である。例えば、先進的な取り組みで知られるSKIP特許事務所では、情報保護を最優先事項とし、Googleの「Gemini Advanced」や「NotebookLM Plus」を導入している<sup>15</sup>。これらのサービスは、ISO 42001等の国際的なセキュリティ規格に準拠し、ユーザーの入力データがモデルのトレーニングやターゲット広告に利用されないことを確約している<sup>15</sup>。また、同事務所は万が一の事態に備えて、賠償責任3億円、サイバーセキュリティ事故対応費用5,000万円をカバーするサイバーリスク補償特約(最高タイプSG)の保険にも加入するなど、極めて慎重な防御策を講じている<sup>15</sup>。しかしながら、このような厳格な環境構築と運用は、すべての事務所や企業が容易に模倣できるものではない。

第二のアプローチは、データを完全に社内に留める「オンプレミス型ローカルLLM」の導入である。日本のITベンダーもこの領域に注力しており、例えばNTTデータは2026年5月より、高セキュリティな閉域網やオンプレミスで動作する知財文書作成エージェントとして「tsuzumi 2」の提供を開始した<sup>16</sup>。また、NECは独自AI「cotomi」を活用した知財DX事業を展開し、特許調査業務を最大94%効率化(22時間の作業を3時間に短縮)するという劇的な成果を挙げている<sup>16</sup>。さらに、知的財産に特化したソフ

トウェアを提供するZanus AIは、15以上の知財モジュールを組み込んだ完全なプライベートAIシステム(Zanus AIオペレーティングシステム)を提供し、先行技術調査や明細書作成をオンプレミス環境で完結させている<sup>20</sup>。

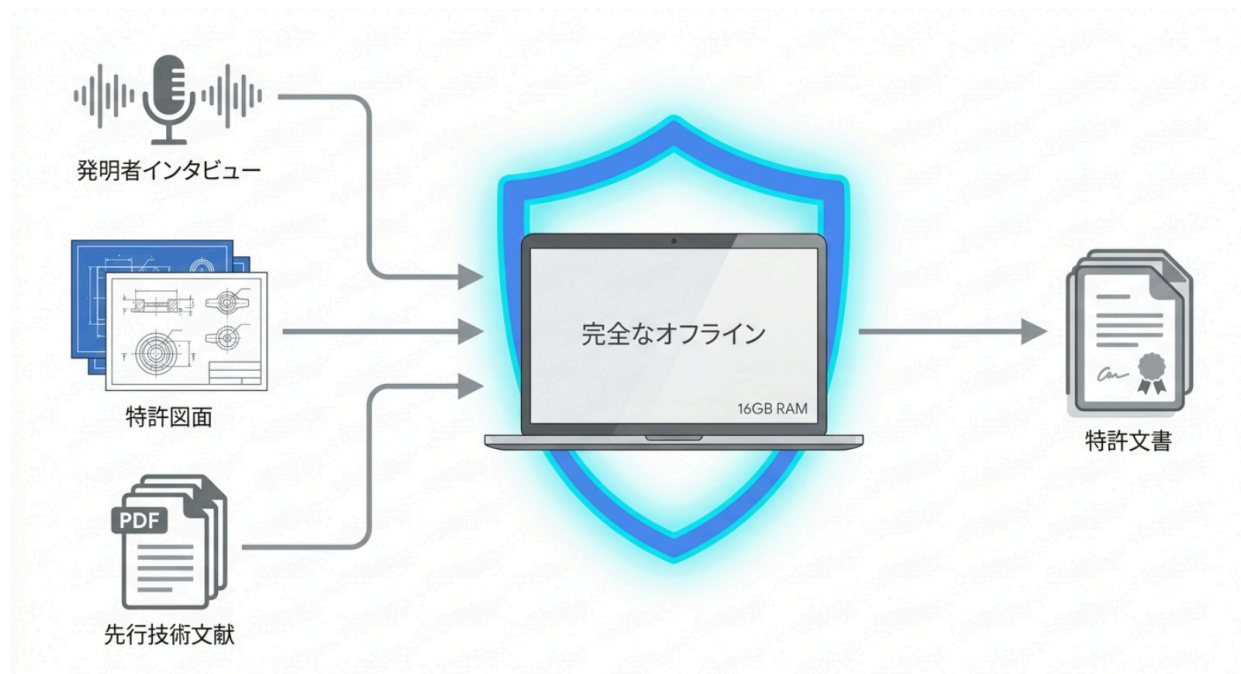
しかし、これらの高性能なオンプレミスLLMソリューションには決定的なデメリットが存在する。それは、高額な初期投資(GPUサーバーの購入、インフラ構築費等)と、保守運用に関する高度な専門知識の要求である<sup>21</sup>。結果として、導入の恩恵は一部の大企業や大規模特許事務所に限定されがちであった。

### 3.2. 日本弁理士会ガイドラインとデータ保護の要請

日本における弁理士の職能団体である日本弁理士会(JPAA)は、生成AIの業務利用に関するガイドライン(弁理士業務AI利活用ガイドライン)を策定・公表している<sup>2</sup>。同ガイドラインでは、弁理士法第75条等との関係を整理しつつ、AIツールを利用する際の利用規約の確認(特に商用利用の可否と学習利用の有無)を厳格に求めている。また、一部の専門家からは、「オプトアウト設定(再学習からの除外)をすれば安全である」「クラウドのサーバーにデータが残ること自体が直ちに新規性喪失につながるわけではない」という冷静な法解釈も提示されているが<sup>2</sup>、実務の現場では、クライアント企業(特に厳格なNDAを要求する大手メーカー)からの「外部クラウドへのデータ送信は一切禁止」というポリシーに従わざるを得ないケースが多い。

Gemma 4 12Bは、この複雑な状況に終止符を打つ可能性を秘めている。なぜなら、追加のインフラ投資を行うことなく、既存の業務用の標準的なノートPC上で、推論やマルチモーダル処理を「完全にネットワークから切断された状態(オフライン環境)」で実行できるからである。通信が発生しない以上、第三者への情報漏洩や再学習のリスクは物理的にゼロとなる。

## Gemma 4 12Bによる完全オフライン・マルチモーダル知財ワークフロー



Gemma 4 12Bは、特許事務所や企業の標準的なノートPC（16GB RAM）上で完全にオフラインで稼働する。発明者のインタビュー音声、複雑な特許図面、膨大な先行技術文献をエンコーダを介さず直接処理することで、最高レベルの機密保持（新規性喪失リスクの排除）を実現する。

### 4. Gemma 4 12Bの技術的特長が知財実務に与える直接的影響

ローカルで安全に稼働するという基本要件を満たした上で、Gemma 4 12Bの卓越したパフォーマンスは、知財業務の各フェーズにおいて、これまで人間が膨大な時間をかけて手作業で行っていたプロセスを根本から変革する。

#### 4.1. 256Kトークンの超長文脈処理による複雑な包袋分析

Gemma 4 12Bは、最大25万6千（256K）トークンという長大なコンテキストウィンドウを備えている（E2BやE4Bは128K）<sup>7</sup>。1トークンを平均して約0.7～1文字の日本語と換算しても、20万文字近くのテキストを一度に処理できる計算となる。さらに、長文脈の理解力を測る「MRCR v2 8 needle 128k」ベンチマークにおいても、12Bモデルと同等以上のモデルが高い精度を示している<sup>8</sup>。

特許業務において「文脈の長さ」は極めて重要である。一つの出願案件を審査段階で処理する際、担当弁理士は「自社の出願明細書（数万文字）」、「特許庁からの拒絶理由通知書」、「引用された複数の先行技術文献（引例1、引例2等の特許公報全文）」、さらには「過去の上申書や補正書の履歴（審査包袋データ）」を同時に比較・検討しなければならない。従来の数千トークン規模しか扱えない

ローカルLLMでは、これらの文書群を一度に入力できず、文書を分割して処理することで文脈のつながりが失われ、精度の低下や幻覚（ハルシネーション）の増加を招いていた。

しかし、Gemma 4 12Bの256Kコンテキストウィンドウを活用すれば、関連するすべての文書や数十件の特許ファミリーのデータ、さらには大規模なソースコードのリポジトリを単一のプロンプトでシームレスに入力することが可能となる<sup>23</sup>。これにより、完全にクローズドなローカル環境において、特定の製品が他社の複数の特許群に抵触しないかを確認する包括的な侵害予防調査（FTO：Freedom To Operate）や、自社特許と競合製品の構成要素の精緻な対比表（クレームチャート）の自動作成が、極めて高い精度で実現する。

## 4.2. 推論モード（Thinking Mode）による拒絶理由通知への高度な対応

特許の実務は、単なるテキストの要約や一般的な文章の生成ではなく、法解釈に基づく厳密な「論理的推論」の連続である。例えば、「引例Aには構成要件Xが開示されているか」「引例Aに引例Bを組み合わせる動機付け（Motivation to combine）が存在し、本願発明は当業者が容易に想到できたか（進歩性欠如）」という論点に対して、隙のない法的なロジックを組み立てる必要がある。

Gemma 4ファミリーは、システムプロンプトの先頭に <think> トークンを含めることで、モデルが最終的な回答を生成する前に、内部でステップ・バイ・ステップの論理的推論プロセスを実行する「思考モード（Thinking Mode / Reasoning Mode）」を搭載している<sup>11</sup>。この機能を有効にすると、モデルはまず <channel>thought\n [内部の推論プロセス] <channel> というタグで囲まれた思考ログを出力し、その後最終的な結論や回答を出力する<sup>11</sup>。

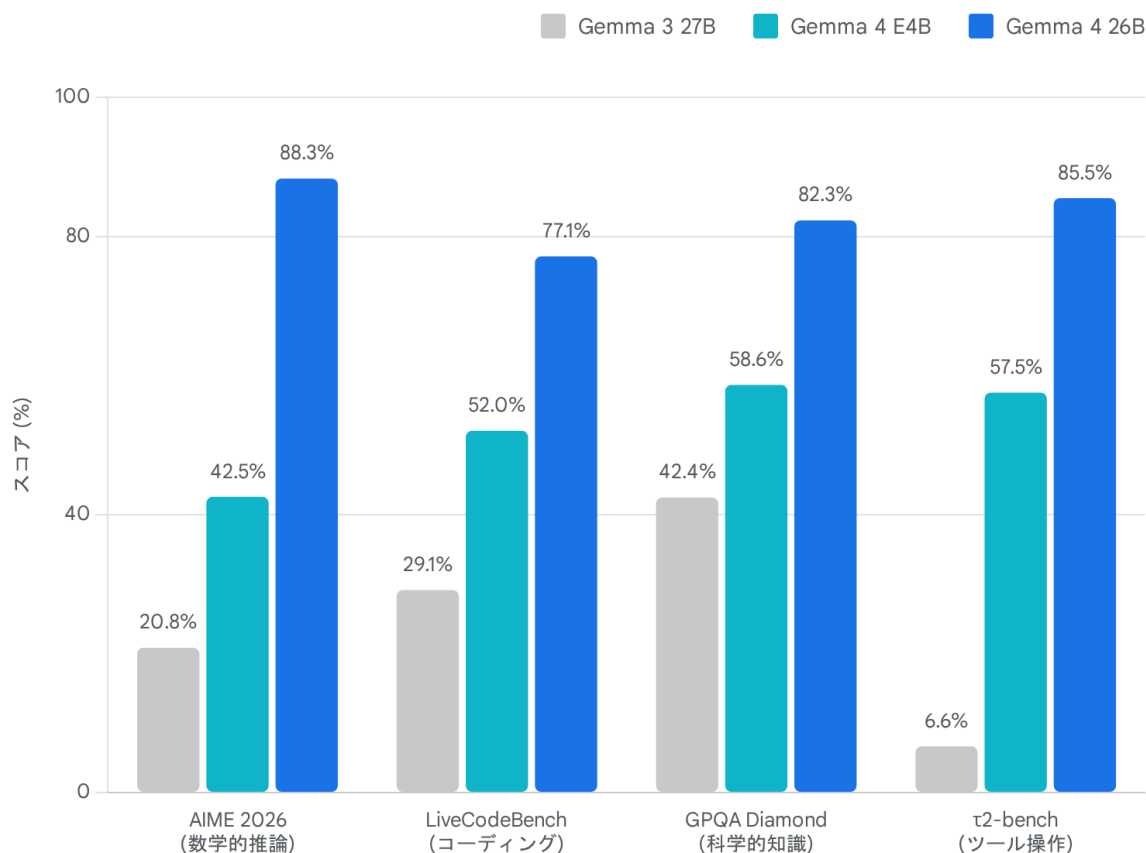
この推論能力は、数学や科学、コーディングの標準的なベンチマークにおいて、旧世代のモデルを圧倒し、より大規模なパラメータを持つ先行モデルに迫る性能を示している。以下の表は、各モデルの主要なベンチマークスコアの比較である<sup>10</sup>。

ベンチマーク指標	テストの性質	Gemma 3 27B IT	Gemma 4 E4B IT Thinking	Gemma 4 26B A4B IT Thinking	Gemma 4 31B IT Thinking
AIME 2026	数学的推論（ツールなし）	20.8%	42.5%	88.3%	89.2%
LiveCodeBench v6	競技プログラミング・コーディング	29.1%	52.0%	77.1%	80.0%
GPQA Diamond	大学院レベルの科学的知識	42.4%	58.6%	82.3%	84.3%
τ2-bench	エージェント的ツール操	6.6%	57.5%	85.5%	86.4%

	作 (Retail)				
--	------------	--	--	--	--

表が示す通り、複雑な論理構築を要求される「数学的推論」や「科学的知識」、そして後述する「ツール操作」において、Gemma 4のThinkingモードは驚異的なスコアを叩き出している。12Bモデルは、メモリフットプリントが26B MoEモデルの半分以下でありながら、これらのベンチマークにおいて26Bモデルに肉薄する性能を発揮すると公式に報告されている<sup>3</sup>。

# Gemma 4ファミリーの推論およびエージェント機能の飛躍的向上



Gemma 4モデル（26B A4Bおよびエッジ向けE4B）と旧世代のGemma 3（27B）のベンチマーク比較。数学的推論（AIME 2026）、科学的知識（GPQA Diamond）、およびエージェント的ツール操作（τ2-bench）において、Gemma 4ファミリーは圧倒的な性能向上を示しており、これが知財業務における複雑な論理構築を可能にする基盤となっている。（注：12Bモデルは26B MoEモデルに迫る性能を発揮する）

データソース: [DeepMind](#), [Google Blog](#)

知財業務において、このThinking Modelは絶大な威力を発揮する。例えば、特許庁からの拒絶理由通知に対する意見書(反論)を作成する場面を想定する。弁理士は、本願明細書、拒絶理由通知書、引用文献A、および引用文献BのテキストをGemma 4 12Bに入力し、「審査官が進歩性を否定した論理構成を分析し、引用文献AとBを組み合わせる動機付けの欠如を証明する法的な反論ロジックを構成せよ」とプロンプトを与える。Gemma 4は、内部の推論プロセス(<|think|>)の中で、「文献Aの課題と文献Bの課題は本質的に異なるか」「本願発明が奏する特有の効果は、AとBの単なる寄せ集めから予測可能か」といった要素を段階的に比較検討する。この推論過程がログとして可視化さ

れるため、弁理士は「AIがなぜその結論に至ったのか」という根拠のトラッキングを容易に行うことができ、弁理士会ガイドラインで求められる「人間による責任あるファクトチェック」の負担が劇的に軽減される<sup>22</sup>。結果として、質の高い意見書のドラフトを短時間で作成することが可能になる。

### 4.3. エージェント機能 (Function Calling) と IDE 連携による明細書作成の自動化

Gemma 4モデルは、単にテキストを生成するだけでなく、外部の構造化されたツールやAPIを呼び出す「Function Calling (関数呼び出し)」をネイティブにサポートしており、自律的なエージェント的ワークフローの構築に最適化されている<sup>8</sup>。

GoogleはGemma 4 12Bのローンチに合わせて、デスクトップ環境でのゼロレイテンシなローカルAI実行を可能にする「LiteRT-LM」を活用したオンデバイス向け開発者統合機能を導入した<sup>4</sup>。この統合により、MacOSアプリ (Google AI Edge Eloquentなど) やローカルのAPIサーバー (litert-lm serve コマンド) としてGemma 4 12Bを簡単に起動できる<sup>4</sup>。

これを特許業務に応用すれば、弁理士は自社の環境内に蓄積された過去の優秀な明細書データベースや辞書データ (いわゆる社内RAGシステム) とGemma 4 12Bを連携させることができる。例えば、Continue、Aider、OpenCodeといったプログラミング用のIDE (統合開発環境) のプラグインにGemma 4 12Bを接続し、ドキュメント作成環境を構築する<sup>4</sup>。弁理士がコアとなるクレーム (特許請求の範囲) を記述すると、AIは社内の過去のボイラープレート (定型文) や独自の表現作法をFunction Calling経由で検索・取得し、各クレームに対応する「発明の詳細な説明 (実施形態)」を段階的かつ自動的に生成していく。特許明細書は時に数十ページから数百ページに及ぶ定型的な論理構造を持つため、このようなエージェントベースの自動化は、作成にかかるリードタイムとコストを抜本的に削減する。

### 4.4. ネイティブ音声・画像処理による発明発掘と図面解釈のシームレス化

特許出願の最初期のフェーズである「発明発掘」では、弁理士が技術者や研究者に対してインタビューを行い、彼らの暗黙知の中から特許化可能な発明のポイント (従来技術の課題、解決手段、顕著な効果) を抽出する作業が行われる。

これまで、このインタビューの録音データをクラウドの音声認識サービス (ASR) に送信することは、営業秘密の保護という観点から強く制限されることが多かった。しかし、Gemma 4ファミリー (特にE2B、E4B、およびマルチモーダル対応モデル) は、140以上の言語で事前学習されており、生の音声信号を直接解釈して自動音声認識および翻訳を行う機能をネイティブに備えている<sup>8</sup>。弁理士がGemma 4 12Bを搭載したノートPCを会議室に持ち込めば、ネットワークから完全に遮断された状態で、技術者の生音声をリアルタイムで文字起こしし、さらにその内容を即座に分析して「発明提案書」のドラフトとして構造化することが可能になる。

さらに、画像処理に関しても、Gemma 4 12Bは対象特許の図面データとテキストを相互に参照し、OCRや多言語の文字認識、UIやチャートの理解を高い精度で実行する<sup>8</sup>。エンジニアが作成した新製品の手書きの設計図やPDFの仕様書をローカルで読み込ませ、対象特許との抵触リスクをスコアリングするといった高度な分析も、外部へのデータ送信なしに実行できる。また、Multi-Token Prediction (MTP) ドラフターが装備されているため、これらの複雑な処理における推論レイテンシも低減されており、ストレスのない対話的プロセスが実現する<sup>3</sup>。

## 5. 競合モデル・代替アーキテクチャとの市場競争力と比較分

# 析

知財向けAIソリューション市場は急速に拡大しており、多様なアプローチが混在している。その中で、ローカル稼働を前提としたGemma 4 12Bはどのようなポジションを占め、どのような優位性を持つのかを比較分析する。

## 5.1. 大規模クラウドモデル(Gemini Advanced等)との棲み分け

Google自身が提供する「Gemini Advanced」や「NotebookLM Plus」などのエンタープライズ向けクラウドサービスは、極めて高い知能と膨大なインデックスデータ(Google PatentsやGoogle Scholarとの連携)を武器としている<sup>15</sup>。これらのクラウドモデルは、インターネット上の最新の論文や海外のデータベースを瞬時に検索し、分析的な推論を統合する能力において、ローカルモデルを凌駕する。SKIP特許事務所のような先進的な組織は、厳しいデータ保護ポリシーの下でこれらのクラウドモデルを積極的に活用し、特許調査や翻訳の能力を「ブースト」している<sup>15</sup>。特に、より高度な推論が可能な「Deep Research」機能や、膨大な社内データを学習させずに参照させるNotebookLMのRAG機能は、強力な武器となる<sup>15</sup>。

したがって、Gemma 4 12Bはこれらのクラウドモデルと競合するものではなく、「相補的な関係」にあると位置づけるべきである。広範な外部データの検索や、社内の大規模なナレッジベース全体の統合にはクラウドモデル(または後述の大型オンプレミス)を使用し、個々の弁理士の手元における極秘情報(新製品の設計図面や発明の着想メモ)の処理、あるいは機密性の高いクライアント先でのオフライン作業にはGemma 4 12Bを使用する、といった「ハイブリッド運用」が今後の知財AI戦略の主流となるだろう<sup>21</sup>。

## 5.2. エンタープライズ向けオンプレミスLLM(tsuzumi, cotomi等)に対する優位性

日本のITベンダーが提供する「tsuzumi 2」(NTTデータ)や「cotomi」(NEC)、あるいは「Zanus AI」のプライベートAIサーバーなどのオンプレミスソリューションは、組織単位でデータを完全に外部から隔離しつつ、大規模なモデルを安定して稼働させることができるというメリットがある<sup>16</sup>。

しかし、これらのソリューションの最大の障壁は「コスト」と「導入のリードタイム」である。オンプレミスで高度なLLMを稼働させるには、数千万円単位の初期投資(GPUサーバーの構築など)と、専任のエンジニアによる保守が必要となる<sup>21</sup>。このため、中堅・中小規模の特許事務所や、予算が限られた企業の知財部では導入が困難であった。

Gemma 4 12Bの圧倒的な優位性は、この「コストの壁」を完全に破壊する点にある。モデル自体はApache 2.0ライセンスで無償公開されており、稼働に必要なハードウェアは「16GBのRAMを積んだ一般的な業務用のノートPC」のみである<sup>3</sup>。追加のハードウェア投資やトークンごとの従量課金、月額サブスクリプション費用が一切発生しない<sup>20</sup>。これは、個人の特許技術者やフリーランスの翻訳者であっても、直ちに世界最高峰のマルチモーダル推論AIを機密環境下で活用できることを意味し、知財業務の生産性向上におけるロングテール市場を一気に開拓する可能性を秘めている。

## 6. 実務導入におけるリスク管理とガバナンス

Gemma 4 12Bがいかにセキュアなローカル環境で動作し、高い推論能力を持つとはいえ、生成AIそのものが内包する本質的なリスクが消滅したわけではない。知財実務家は、専門家としての高度な倫理観と厳格なガバナンスをもって、このテクノロジーを運用する責任がある。

## 6.1. ハルシネーションに対するファクトチェックと専門家の責任

AIモデルは、最もらしいが事実とは異なる情報(ハルシネーション)を生成する性質を持っている。SKIP特許事務所の検証事例でも明確に指摘されているように、高度なAIであっても、特許調査において存在しない文献番号をでっち上げたり、架空の根拠となるURLを生成したりする事象が確認されている<sup>15</sup>。

日本弁理士会が公表したガイドラインにおいても、情報の正確性に関する警告がなされており、「情報が古い場合もあり得るため、アップデートされた情報に基づいたファクトチェックを行うことも重要となる」「生成結果を利用するにあたっては、弁理士はその内容について責任をもって検討・確認をし、判断しなければならない」と明記されている<sup>22</sup>。

Gemma 4 12Bの「Thinking Mode」は、推論プロセスを可視化することでこのファクトチェックを支援するが、生成されたクレームの文言が本当に先行技術を回避して広い権利範囲を確保できているか、あるいは反論のロジックが特許審査基準に照らして妥当であるかについては、最終的に人間(弁理士)の専門的判断が不可欠である。AIはあくまで「極めて高速で有能な下書き作成者(ドラフター)」として位置づけ、法的見解を提供する主体は人間でなければならない<sup>32</sup>。

## 6.2. 著作権等の第三者権利への配慮と安全な運用環境の構築

同ガイドラインでは、生成物の利用に関して「他者の知的財産権(特に著作権)を侵害していないか確認が必要である」とも指摘している<sup>22</sup>。特許明細書の文章自体に著作権の懸念が生じるケースは実務上稀であるが、AIがRAGなどで参照した他社の論文や文献の独自表現をそのまま出力し、それを商業的に利用した場合には法的問題に発展する可能性がある。Gemma 4 12Bのライセンス自体は寛容なApache 2.0であるが<sup>3</sup>、生成物に対する最終的な法的責任はユーザーが負うという大原則を忘れてはならない。

さらに、「ノートPCでのローカル稼働」とはいえ、そのPC自体がインターネットに常時接続され、マルウェアに感染しているような状況であれば、Gemma 4 12Bの処理内容がバックドアを通じて外部に漏洩するリスクがある<sup>21</sup>。真に機密性の高い未出願の発明を扱う場合は、Gemma 4 12Bを実行する専用の端末をネットワークから物理的または論理的に切断する(エアギャップの構築)ことや、組織的なEDR(Endpoint Detection and Response)等の高度なセキュリティ監視を徹底することが、エンタープライズ水準の運用においては必須となる。

## 7. 結論: 知財実務家の役割の再定義と高付加価値化への道

Gemma 4 12Bの登場は、知的財産業務における生成AIの活用において、長らく業界を悩ませてきた「セキュリティの確保」と「導入・運用コストの低減」という二律背反の課題を同時に解決する、極めて重要な技術的ブレイクスルーである。

16GBのRAMを備えた標準的なノートPC上で稼働するコンパクトさをもちながら、エンコーダフリーアーキテクチャによるシームレスな画像(特許図面)・音声(発明ヒアリング)の直接処理、256Kトークンという長大なコンテキストウィンドウによる包括的な包袋・文献分析、そしてThinking Modelによる深遠な論理的推論能力を併せ持つこのモデルは、特許事務所や企業の知財部門にとって「完全に機密が保持された、専属の超優秀な特許技術者」を無料で手に入れることに等しい。

音声によるリアルタイムの発明発掘から、機能呼び出し(Function Calling)とIDE連携を活用した明細書作成の自動化、さらには複雑な拒絶理由通知の論理構成を分析して反論を構築するプロセスに至るまで、Gemma 4 12Bは知財ワークフローのあらゆるフェーズを根底から加速・高度化させるポ

テンシヤルを有している。

しかしながら、この強力なツールを実務で安全に運用するためには、日本弁理士会のガイドラインに則した厳格なファクトチェック体制の構築と、AIの生成プロセスを監査する専門家としての高い倫理観が不可欠である。Gemma 4 12Bは決して知財実務家の仕事を奪うものではなく、むしろ定型的な文書作成や情報の検索といった労働集約的な作業から専門家を解放するものである。これにより、弁理士や知財担当者は、より高度な特許ポートフォリオ戦略の立案や、クライアントの事業価値を最大化するための知財コンサルティングといった、人間にしかできない高次元の意思決定と価値創造の領域へと、その役割を再定義していくことになるだろう。

## 引用文献

1. 生成AIを特許業務に活用する難しさとは？ | 平井智之/リーガルテックCEO - note, 6月7, 2026にアクセス、[https://note.com/yutori\\_jd/n/n3831a1e6769b](https://note.com/yutori_jd/n/n3831a1e6769b)
2. 第4回 ChatGPTの利用は新規性を喪失するのか？ | WEB会議で全国 ..., 6月7, 2026にアクセス、<https://takayama-patent.com/archives/2723>
3. Introducing Gemma 4 12B - Google Blog, 6月7, 2026にアクセス、<https://blog.google/innovation-and-ai/technology/developers-tools/introducing-gemma-4-12b/>
4. Gemma 4 12B: The Developer Guide, 6月7, 2026にアクセス、<https://developers.googleblog.com/gemma-4-12b-the-developer-guide/>
5. Gemma releases | Google AI for Developers, 6月7, 2026にアクセス、<https://ai.google.dev/gemma/docs/releases>
6. Google's latest on-device AI model is custom-made for your laptop, 6月7, 2026にアクセス、<https://www.androidauthority.com/google-gemma-4-12b-multimodal-ai-model-3674379/>
7. Google's new open source Gemma 4 12B analyzes audio, video — and runs entirely locally on a typical 16GB enterprise laptop, 6月7, 2026にアクセス、<https://venturebeat.com/technology/googles-new-open-source-gemma-4-12b-analyzes-audio-video-and-runs-entirely-locally-on-a-typical-16gb-enterprise-laptop>
8. Gemma 4 - LM Studio, 6月7, 2026にアクセス、<https://lmstudio.ai/models/gemma-4>
9. google/gemma-4-12B - Hugging Face, 6月7, 2026にアクセス、<https://huggingface.co/google/gemma-4-12B>
10. Gemma 4 - Google DeepMind, 6月7, 2026にアクセス、<https://deepmind.google/models/gemma/gemma-4/>
11. gemma4:12b - Ollama, 6月7, 2026にアクセス、<https://ollama.com/library/gemma4:12b>
12. gemma4 - Ollama, 6月7, 2026にアクセス、<https://ollama.com/library/gemma4>
13. Google Gemma4 12B released, 6月7, 2026にアクセス、<https://medium.com/data-science-in-your-pocket/google-gemma4-12b-released-b497e9c063e4>
14. Gemma 4 QAT models: Optimizing model compression for mobile and laptop efficiency - Google Blog, 6月7, 2026にアクセス、

- <https://blog.google/innovation-and-ai/technology/developers-tools/quantization-aware-training-gemma-4/>
15. 【弁理士業務AI利活用ガイドライン準拠！】SKIPの全社員がGoogle ..., 6月7, 2026にアクセス、<https://skiplaw.jp/blog/13444/>
  16. 2026年知財分野における日本の国産LLM採用動向レポート, 6月7, 2026にアクセス、  
<https://yorozuipsc.com/uploads/1/3/2/5/132566344/d1d63f2a696b2774dbd5.pdf>
  17. Microsoft Azure上で「tsuzumi 2」の提供開始、生成AIの検証環境を整備 | NTTデータグループ, 6月7, 2026にアクセス、  
<https://www.nttdata.com/global/ja/news/topics/2026/052001/>
  18. NTT版LLM tsuzumi 2アップデート ~世界トップレベルの図表入り日本語ビジネス文書処理性能を1GPU環境で実現~ | ニュースリリース, 6月7, 2026にアクセス、  
<https://group.ntt.jp/newsrelease/2026/05/19/260519a.html>
  19. NECが知財AI開発で実現した「最大94%効率化」。特許調査は22時間から3時間へ, 6月7, 2026にアクセス、  
<https://www.businessinsider.jp/article/2601-nec-ai-intellectual-property-efficiency/>
  20. 知的財産・特許事務所向けAIソフトウェア | オンプレミス型AI - Zanus AI, 6月7, 2026にアクセス、  
<https://ja.zanusai.com/products/ai-software-for-ip-patent>
  21. ローカルLLM・オンプレミスLLMとは？企業が“AIを内側に置く”理由 | DXPOカレッジ, 6月7, 2026にアクセス、  
<https://dxpo.jp/college/ai-agent/local-llm.html>
  22. 弁理士業務AI利活用ガイドライン, 6月7, 2026にアクセス、  
<https://www.jpaa.or.jp/cms/wp-content/uploads/2025/04/AIservices-guideline.pdf>
  23. Gemma 4: Byte for byte, the most capable open models - Google Blog, 6月7, 2026にアクセス、  
<https://blog.google/innovation-and-ai/technology/developers-tools/gemma-4/>
  24. Thinking mode in Gemma | Google AI for Developers, 6月7, 2026にアクセス、  
<https://ai.google.dev/gemma/docs/capabilities/thinking>
  25. google/gemma-4-31B-it - Hugging Face, 6月7, 2026にアクセス、  
<https://huggingface.co/google/gemma-4-31B-it>
  26. Found how to toggle reasoning mode for Gemma in LM-Studio! : r/LocalLLaMA - Reddit, 6月7, 2026にアクセス、  
[https://www.reddit.com/r/LocalLLaMA/comments/1sc9ucc/found\\_how\\_to\\_toggle\\_reasoning\\_mode\\_for\\_gemma\\_in/](https://www.reddit.com/r/LocalLLaMA/comments/1sc9ucc/found_how_to_toggle_reasoning_mode_for_gemma_in/)
  27. Gemma 4 12B: Specs, Benchmarks & How to Run It Locally - Build Fast with AI, 6月7, 2026にアクセス、  
<https://www.buildfastwithai.com/blogs/gemma-4-12b-guide>
  28. How to Enable Gemma 4 Thinking Mode in LM Studio and OpenCode - Antonio Leiva, 6月7, 2026にアクセス、  
<https://antonioleiva.com/enable-gemma-thinking-mode-lm-studio-opencode>
  29. NEW Google Gemma 4 12B AI Update 🤖, 6月7, 2026にアクセス、  
<https://www.youtube.com/watch?v=mZeCrrbxX28>
  30. Gemma 4 12B - Google's Unified Multimodal Model Running Locally, 6月7, 2026にアクセス、  
<https://www.youtube.com/watch?v=Uh08fdDjlpU>
  31. オンプレミス生成AI基盤 — データを外に出さないローカルLLM環境 - 製品詳細 - 2026年5月大阪開催, 6月7, 2026にアクセス、

<https://www.nepconjapan.jp/osaka/ja-jp/search/2026/product/product-details.exh-fe14fbfe-109a-41e0-a0f3-031b6bb83a59.%E3%82%AA%E3%83%B3%E3%83%97%E3%83%AC%E3%83%9F%E3%82%B9%E7%94%9F%E6%88%90ai%E5%9F%BA%E7%9B%A4%20%20%E3%83%87%E3%83%BC%E3%82%BF%E3%82%92%E5%A4%96%E3%81%AB%E5%87%BA%E3%81%95%E3%81%AA%E3%81%84%E3%83%AD%E3%83%BC%E3%82%AB%E3%83%ABIm%E7%92%B0%E5%A2%83.pro-46dcf51e-cd4a-4d5d-94ec-b1d4aff4ecbe.html>

32. 生成AI導入のお知らせ - IPTech弁理士法人, 6月 7, 2026にアクセス、  
<https://iptech.jp/info/250328>