

# ARC-AGI-2の成績から読み解く知的財産業務における大規模言語モデルの性能差異と実務への影響

Gemini 3.1 pro

## 1. 序論：流動性知能と特許実務の交差点

人工知能(AI)、とりわけ大規模言語モデル(LLM)の進化は、自然言語処理の枠を超え、高度な専門知識と論理的推論が要求される知的財産(IP)実務の領域に根本的なパラダイムシフトをもたらしている。歴史的に、特許明細書の作成、先行技術調査、および無効審判における緻密な論理構築などの知的財産業務は、技術的理解と厳密な法的制約の双方を統合する「人間の高度な認知能力」の独壇場であった。日本弁理士会(JPAA)の報告にも示されるように、日本の労働人口の減少、国立大学法人化に伴う基礎研究能力の相対的な変化、および若手弁理士の採用難というマクロ環境の圧力の中、限られたリソースで権利価値を最大化するためのAI技術の導入は喫緊の課題となっている<sup>1</sup>。しかしながら、最新の推論型AIモデルの登場と旧世代モデルの混在により、知財実務におけるAIの適用限界と可能性の境界線は複雑化している。

このAIの推論能力を客観的かつ厳密に測定する試みとして、François Cholletらによって提唱された「ARC-AGI(Abstraction and Reasoning Corpus for Artificial General Intelligence)」が極めて重要な指標として機能している<sup>5</sup>。特に2025年にリリースされた「ARC-AGI-2」は、事前の暗記や計算機的なブルートフォース(総当たり)探索が通用しないように意図的に設計されており、未知の問題に対する適応力、すなわち「流動性知能(Fluid Intelligence)」を直接的に評価するベンチマークである<sup>6</sup>。流動性知能とは、過去の経験や蓄積された知識(結晶性知能)に依存せず、新規な問題に対してその場でパターンを見出し、論理を構築して解決する能力を指す<sup>8</sup>。

本稿では、ARC-AGI-2のスコアにおいて「高性能(高スコア)を示す推論特化型LLM」と、「低性能(低スコア)に留まる従来の汎用ベースLLM」との間に、実際の知的財産業務(特に特許実務)においてどのような質的・構造的な違いが表れるのかを網羅的に分析する。ARC-AGI-2が測定する「シンボリック解釈(Symbolic Interpretation)」や「合成推論(Compositional Reasoning)」といった能力は、特許請求の範囲(クレーム)の解釈や、新規性・進歩性の判断における論理構築能力と深く関連している<sup>8</sup>。この相関を解き明かすことで、知財実務におけるAI導入の現在地と、内在するハルシネーションなどのリスク、そして将来の費用対効果(コスト・オブ・インテリジェンス)の最適化戦略を提示する。

## 2. ARC-AGI-2のアーキテクチャと流動性知能の測定基準

知的財産業務におけるLLMの性能差異を理解するためには、まずARC-AGI-2がいかなる認知能力を測定しているかを詳まびらかにする必要がある。ARC-AGI-2は、汎用人工知能(AGI)に向けた進捗を測定するため、従来型の知識集約型ベンチマーク(MMLUなど)とは全く異なる設計思想を持つ<sup>7</sup>

。

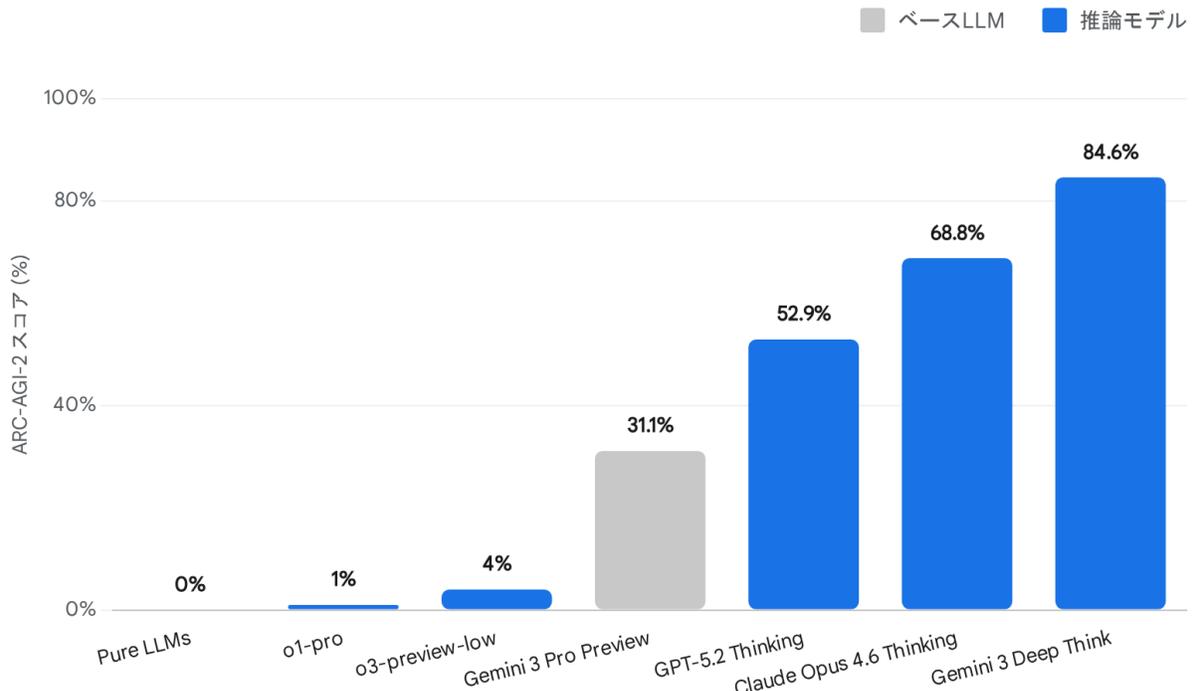
## 2.1 ベンチマークの構造と人間のベースライン

ARC-AGI-2のテストは、1x1から最大30x30のグリッド上に、最大10色のカラーブロックが配置された視覚的なパズルで構成されている<sup>6</sup>。AIシステムは、少数(通常2~5個)の入力と出力のデモンストレーション画像から、その背後に隠された抽象的な変換規則(潜在的マッピング)を推論し、それを一度も見ることがない新しいテスト入力に適用して、ピクセル単位で完全に一致する正解グリッドを生成しなければならない<sup>5</sup>。これは極端な少数ショット(Few-shot)の汎化能力を要求する<sup>5</sup>。

このベンチマークは、意図的に「人間にとっては容易であるが、AIにとっては困難」に設計されている<sup>12</sup>。開発陣による厳格な校正プロセスを経ており、事前の人間によるテストでは、専門知識を持たない400人のボランティアが1,417のユニークなタスクに挑戦し、すべてのタスクにおいて「少なくとも2人の人間が2回以内の試行で正解する」ことが確認されている<sup>12</sup>。人間の平均解答時間は1タスクあたり約2.3分であり、全体の正答率は100%に達する<sup>12</sup>。また、解答者の年齢、性別、学歴、専門的経験などの人口統計学的要因とパフォーマンスとの間に相関関係は見られず、これが特定の専門知識ではなく、普遍的な抽象推論能力を測定していることを証明している<sup>12</sup>。

# 流動性知能（ARC-AGI-2）と高度知財タスク対応能力の 相関

主要AIモデルの ARC-AGI-2 評価スコア比較



ARC-AGI-2における抽象的推論スコアの飛躍は、単なるパズル解決能力ではなく、特許実務など高度な論理的推論を要する専門領域における実務遂行能力の向上と強く連動している。

データソース: [ARC Prize Foundation](#), [Intuition Labs](#), [Implicator.ai](#)

## 2.2 純粋な言語モデルの限界と推論モデルの飛躍

ARC-AGI-2におけるAIシステムのパフォーマンスは、システムのアーキテクチャによって劇的に二極化している。事前学習によるパターンの暗記に大きく依存する純粋な大規模言語モデル(GPT-4、Claude 3.5 Sonnetなど)は、ARC-AGI-2の厳密な環境下ではスコアが0%~5%未満にとどまり、事実上の完全な破綻を示している<sup>8</sup>。これは、データセット内のすべてのタスクがユニークであり、過去の学習データからの表面的なテキスト検索やパターンマッチングでは解決できないためである<sup>6</sup>。

一方で、テスト時計算(Test-Time Compute)を導入し、思考の連鎖(Chain-of-Thought: CoT)や強化学習(RL)によって内部で推論ツリーを探索する最新の「推論特化型モデル(Reasoning Models)」は、顕著な飛躍を遂げている<sup>6</sup>。例えば、OpenAIのGPT-5.2はARC-AGI-2において52.9%~54%のスコアを記録し<sup>12</sup>、さらにGoogleのGemini 3 Deep Thinkは84.6%という驚異的な精度を独立検証

で達成した<sup>18</sup>。また、強化学習のみで訓練されたDeepSeek R1-Zeroや、教師あり学習を組み合わせたDeepSeek R1も、オープンソースモデルでありながら15%~20%のスコアを叩き出し、純粋なLLMの限界を突破している<sup>14</sup>。

## 2.3 ブルートフォースへの耐性と効率性の評価

前身であるARC-AGI-1の課題の一つは、Kaggleコンペティションのトップ解答の多くが、知能の証明ではなく、計算リソースに物を言わせたブルートフォース(プログラムの総当たり探索)に依存していたことであった。プライベート評価セットの約49%がこの手法で解かれてしまい、ベンチマークとしてのシグナルが希釈されていた<sup>6</sup>。

ARC-AGI-2はこれを意図的に排除する設計を採用している。タスクの認知複雑性を高め、計算量に依存した素朴な探索では解けないようにしている<sup>6</sup>。さらに重要な点として、ARC-AGI-2は正答率だけでなく、「タスクあたりのコスト(効率性)」を重要な評価基準として導入した<sup>8</sup>。真の知能とは、無尽蔵の計算資源を使って正解を探し当てることではなく、人間の脳のように最小限のエネルギーと認知負荷で効率的に未知のルールを抽出・適応することである、という哲学が背景にある<sup>9</sup>。

## 3. 抽象的推論能力と知的財産業務の認知的交差点

ARC-AGI-2が低性能LLMを脱落させ、高性能LLMのみがクリアできる「中核的な認知的挑戦」は、主に3つのカテゴリに分類される<sup>8</sup>。これらは単なるパズル上のギミックではなく、特許実務において弁理士や審査官が日常的に行っている高度な法務・技術的推論プロセスと直接的な相似形を成している。

### 3.1 シンボリック解釈(Symbolic Interpretation)とクレーム構築

第一の課題は「シンボリック解釈」である。低性能なモデルは、グリッド上のピクセルを単なる視覚的パターン、対称性、あるいは物理的な形状の移動としてしか捉えることができない<sup>8</sup>。しかし、ARC-AGI-2の高度なタスクでは、特定の色や形が、ある時は「移動を示すベクトル」であり、別の時は「論理ゲート(AND/OR)」を意味し、さらに別の時は「状態を切り替えるスイッチ」であるといったように、その要素が持つ「機能的・意味論的(セマンティック)な役割」を文脈から読み取らなければならぬ<sup>12</sup>。高性能LLMは、このようなシンボルに意味を付与し、その意味に基づいて全体を操作することができる<sup>8</sup>。

特許実務において、これはまさに「クレーム解釈(Claim Construction)」そのものである。特許請求の範囲に記載される「係合手段」「付勢部材」「情報処理装置」といった用語は、日常語としての表面的な辞書的意味を持つのではなく、明細書全体や図面、さらには出願経過(ファイルラッパー)というコンテキストの中で定義される「法的な境界線を持つ特有のシンボル」である<sup>21</sup>。低性能LLMは、特許文書を処理する際に、これらの用語を一般的なウェブテキストに基づく統計的な確率分布(パターンマッチング)として処理するため、特許法上の厳密な権利範囲を誤認したり、機能的クレームを構造的制約と混同したりする<sup>21</sup>。一方、高性能LLMは、未知の専門用語であっても、明細書の「詳細な説明」部分からその機能的・構造的意味を動的に推論(流動性知能を発揮)し、一貫した法的シンボル

として維持しながら侵害可否や新規性の判断を行う能力に長けている。

### 3.2 合成推論 (Compositional Reasoning) と特許の階層構造

第二の課題は「合成推論」である。これは、複数の異なるルールを同時に適用し、それらが互いに干渉・相互作用する状況で正しい結果を導き出す能力を指す。ARC-AGI-2の開発レポートによれば、単一のグローバルルール(例:全体を90度回転させる)であれば低性能LLMでも一貫して発見・適用できる。しかし、複数のルール(例:赤いブロックは右に移動するが、青いブロックと衝突した場合は色が混ざって上に移動する)が絡み合うタスクでは、途端に推論が破綻する<sup>8</sup>。

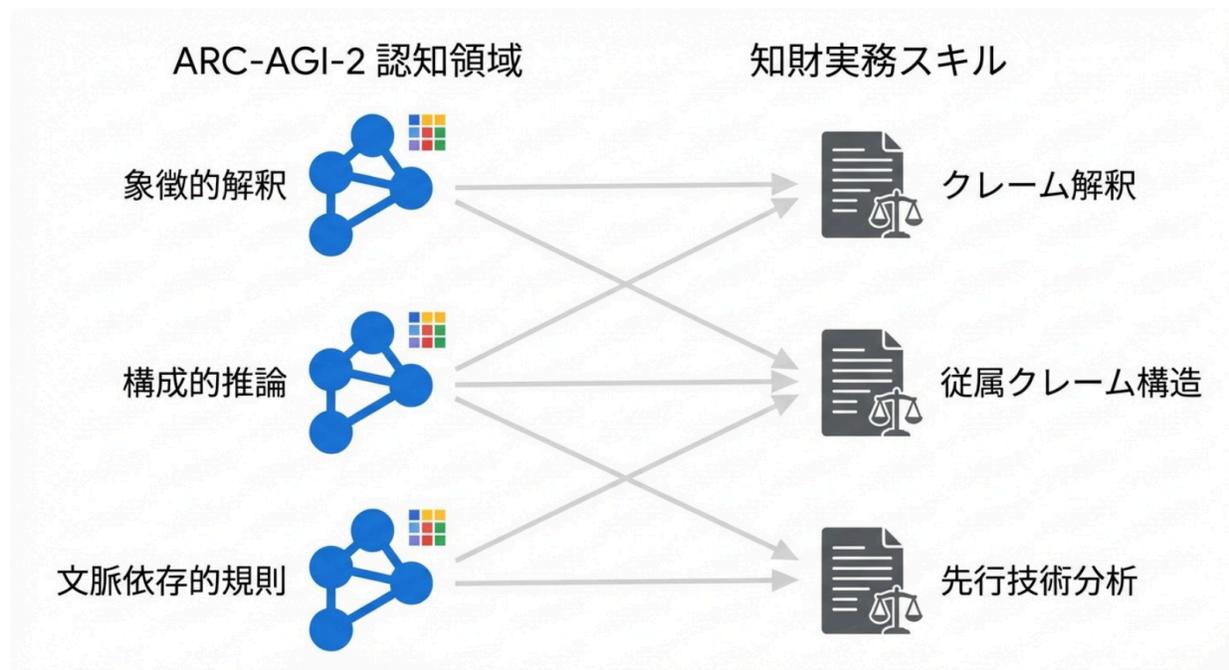
特許明細書、特に特許請求の範囲は、独立クレームに対して多数の従属クレームが連なる、極めて高度な「合成論理構造」のテキストである。従属クレームは、独立クレームのすべての構成要件を継承しつつ、新たな限定的要素(ルール)を追加する。低性能LLMに特許案の作成や分析を行わせると、複数の制約を同時に保持できず、先行する構成要件との整合性(先行詞の要件: Antecedent Basis)を見失ったり、独立クレームの前提と決定的に矛盾する従属クレームを生成したりする致命的なエラーを引き起こす<sup>21</sup>。高性能LLMは、複数の制約条件(技術的整合性と法的要件)を合成的に推論し、全体の論理的完全性(Hierarchy and Dependency)を保持しながらテキストを生成・評価する能力において圧倒的な優位性を持つ。

### 3.3 コンテキスト依存のルール適用と先行技術対比

第三の課題は、コンテキスト依存のルール適用(Contextual Rule Application)である。ARC-AGI-2では、同一の形状やパターンであっても、周囲の状況や全体の配置という「コンテキスト」によって適用すべきルールを動的に切り替える必要がある<sup>8</sup>。低性能モデルは、根底にある選択原理を理解するのではなく、表面的なパターンに固執するため、この種のタスクで失敗する傾向が強い<sup>12</sup>。

この能力の差異は、特許審査における先行技術文献(Prior Art)の対比や進歩性(Non-obviousness)の判断において決定的な違いを生む。先行技術に開示されているある要素(例:特定の合金材料)が、対象発明においても同じ動機付けで適用できるかどうかは、その発明が解決しようとする課題や、技術分野全体のコンテキストに依存する。低性能LLMは、単なるキーワードの類似度(Exact Matchや意味的ベクトル距離)に基づいて「関連性あり」と表面的な判断を下すため、特許審査官の論理とは程遠い結果を出力する<sup>26</sup>。一方、高性能な推論モデルはコンテキストの差異を評価し、異なる分野の技術を組み合わせる際の「阻害要因」や「動機付けの有無」といった、当業者の論理的推論プロセスをシミュレートする能力を備えつつある。

## ARC-AGI-2の認知タスクと特許実務における要求スキルの相関構造



ARC-AGI-2で測定される流動性知能（左側）は、特許明細書の構造的・法的完全性を担保するための実務スキル（右側）と直接的に対応している。低スコアのLLMは表面的な言語生成に留まるのに対し、高スコアのモデルはこれらの深層論理を模倣できる。

### 4. 低性能LLM（ベースモデル）が特許実務にもたらす限界と致命的リスク

ARC-AGI-2において数パーセント未満のスコアしか出せない、いわゆる「純粋な自己回帰型の大規模言語モデル（言語ファーストモデル）」は、日常的な文章作成や要約においては一見して流暢なテキストを生成する。しかし、その内部機構の限界から、知的財産業務の最前線においては深刻なリスクと機能不全を呈する<sup>13</sup>。

#### 4.1 言語的流暢性と論理的破綻の乖離

低性能LLMは、本質的に「これまでの文脈に基づいて、次に続く確率が最も高い単語（トークン）を予測する」統計的パターンマッチングマシンである<sup>7</sup>。これらは特許の技術分野や一般的な法用語に関する豊富な知識（結晶性知能）を事前学習により保持しているため、特許出願のドラフティングにおいて、体裁の整った「特許らしい」文章を高速で生成することは得意である<sup>21</sup>。

しかし、このような「言語ファースト（Language-First）」のツールは、ドラフティングを純粋な執筆作業

(ライティングタスク)としてしか捉えていない<sup>21</sup>。科学的事実の評価や論理の因果関係を自律的に検証する能力を持たないため、もっともらしいが技術的に成立しない構造や、自然法則に反するメカニズムを平然と出力する「ハルシネーション(幻覚)」を引き起こす<sup>29</sup>。特許法において、明細書は当業者がその発明を実施できる程度に明確かつ十分に記載されなければならない(実施可能要件: Enablement Requirement)。LLMが生成した架空の実験データや、動作不可能な実施例をそのまま出願すれば、特許は拒絶されるか、権利化後に容易に無効化されてしまう<sup>22</sup>。

## 4.2 先行詞基準(Antecedent Basis)と構造的欠陥

特許請求の範囲(クレーム)は、技術的特徴の羅列ではなく、各要素の接続関係と法的境界を定義する厳密な階層的テキストである<sup>22</sup>。特に米国特許実務などにおいて厳しく問われる「先行詞基準(Antecedent Basis)」は、ある用語がクレーム内で「The」や「Said(前記)」を伴って再登場する場合、それが以前に明確に定義された同一の要素を指していなければならないという原則である<sup>22</sup>。

PatentScoreやPat-DEVALといったLLM生成特許の多次元評価フレームワークを用いた研究では、低性能LLMを用いた場合、この先行詞基準の欠如や構造的矛盾が多発することが確認されている<sup>22</sup>。例えば、クレーム1で「第1の通信部」と定義したにもかかわらず、従属クレームで「前記送受信部」という異なる用語を突如として使用するなどである。これは、低性能LLMが長文を生成する際、ローカルな注意機構(Attention)に引きずられ、ドキュメント全体のグローバルな論理構造を維持する「合成推論」の能力が不足していることに起因する<sup>21</sup>。このような不明確性(Indefiniteness)は、権利範囲の解釈において致命的な弱点となる<sup>22</sup>。

## 4.3 化学・バイオ分野におけるマーカッシュ構造の崩壊

AIの推論限界が最もクリティカルな結果をもたらすのが、化学およびバイオテクノロジー分野の特許実務である。この分野では、基本骨格に対して複数の置換基のバリエーションを包括的に表現する「マーカッシュ群(Markush Groups)」が多用される<sup>21</sup>。

低性能LLMは、IUPAC命名法や化学式の文字列を自然言語のテキストデータとして確率的に出力する。そのため、化学的推論を伴わない生成結果には、「原子価(Valency)の異常」、「結合位置の矛盾」、「立体化学的に存在し得ない足場構造(Scaffold)」といったエラーが頻発する<sup>21</sup>。弁理士が「この基本骨格に基づき、有効な置換基のバリエーションを展開せよ」と指示した場合、言語ファーストのLLMは無数の組み合わせを生成するが、その中に化学的に成立しない要素や、明細書本文(Detailed Description)に一切のサポート(実施例や薬理データ)がない要素を意図せず紛れ込ませる<sup>21</sup>。これは、サポート要件(Written Description)の違反に直結する。マーカッシュ構造の一部が無効であれば、クレーム全体、あるいは特許出願そのものの法的な健全性が損なわれるため、この領域において低性能LLMをドラフティングの主力とすることは極めて危険である<sup>21</sup>。

## 4.4 IPBenchが示す汎用モデルの分類・評価能力の限界

知的財産タスクに特化した包括的な多言語ベンチマークである「IPBench」の調査結果は、これまでの経験則を定量的に裏付けている。IPBenchは、WebbのDOK(Depth of Knowledge)理論に基づき、知財タスクを「情報処理(Information Processing)」、「論理推論(Logical Reasoning)」、「判別

評価(Discriminant Evaluation)」、「創造的生成(Creative Generation)」の4階層に分類し、特許分類(IPC/CPC)、権利帰属分析、無効性特定など20の個別タスクを設けている<sup>26</sup>。

このベンチマークにおいて、16種類の最先端LLM(GPT-4o、Claude、Llama等のチャットベースの汎用モデルを含む)をテストした結果、最も高性能なモデルであっても全体の正答率は75.8%にとどまった<sup>26</sup>。特に注目すべきは、国際特許分類(IPC)や共通特許分類(CPC)の完全一致(Exact Match)タスクにおいて、低性能モデルが致命的な結果(Disaster)を示したことである<sup>26</sup>。特許分類の付与は、単にタイトルや要約のキーワードを拾うだけではなく、発明の技術的課題、解決手段、および用途を抽象化し、極めて細分化された階層構造を持つ分類体系に正確にマッピングする作業である<sup>26</sup>。低性能LLMは、ARC-AGI-2における「シンボリック解釈」の欠如と同様に、特許文書の技術的本質を適切に抽象化して分類シンボルに落とし込むことができず、表面的な語彙の類似性に騙されて誤った分類を出力してしまうのである<sup>26</sup>。

さらに、IPBenchのエラー分析では、全エラーの33%が「推論エラー(Reasoning Error)」に分類された。法規定の変遷や地域による知財制度の差異(例:米国と中国の特許法の違い)を動的に適用するタスクにおいても、推論能力の欠如が低いパフォーマンスの主要因となっている<sup>32</sup>。

## 5. 高性能LLM(推論特化型モデル)による知財業務のパラダイムシフト

これに対し、ARC-AGI-2において極めて高いスコア(50%~80%台)を記録したモデル群(OpenAI o3、Gemini 3 Deep Think、DeepSeek R1など)は、知的財産業務の質を根本から変革し、単なる「執筆補助ツール」から「高度な論理検証パートナー」へとAIの役割を押し上げるポテンシャルを秘めている<sup>14</sup>。これらのモデルの躍進を支えているのは、「テスト時計算(Test-Time Compute)」、「思考の連鎖(CoT)」、そして「強化学習(RL)」による自己修正能力である<sup>14</sup>。

### 5.1 「Structure-First」アプローチと論理の自己検証(Test-Time Compute)

高性能LLMは、特許明細書を単なる「文章(Language)」としてではなく、「制約を満たすべき論理構造(Structure)」として扱う能力を獲得している(Structure-Firstアプローチ)<sup>21</sup>。

ARC-AGI-2で高スコアを出すための鍵は、推論時(テスト時)に追加の計算資源を投下し、モデル自らが長い思考の連鎖(CoT)を生成することにある<sup>14</sup>。この過程で、モデルは自ら生成した仮説(ルールの推論)を別の入力例に適用して正しさを検証し、間違っていればバックトラック(後戻り)して別の探索経路を試みる<sup>17</sup>。例えば、Sakana AIが開発した「AB-MCTS(Adaptive Branching Monte Carlo Tree Search)」のようなアルゴリズムは、LLMに効果的な試行錯誤を行わせ、解の探索空間を最適化する<sup>34</sup>。

知財実務において、この能力は特許案の自律的な品質保証メカニズムとして機能する。高性能LLMを用いた特許作成システムは、クレームを生成した後に、自らが作成したクレームツリーを逆に解析する。「要素Aが独立クレームで定義されていない」「実施例の図面3に記載されていない特徴が従属クレームで追加されている」「マーカッシュ群の定義が明細書本文の実施例と整合していない」とい

た内部矛盾を自律的に検出し、修正することが可能になる<sup>21</sup>。さらに、SymbolicaのAgenticaフレームワークのような、LLMとPython REPL(対話型評価環境)を結合したニューロシンボリック(Neuro-symbolic)なアプローチは、LLMが自らの推論をコードとして実行し、検証することを可能にする<sup>5</sup>。化学特許の例であれば、LLMが生成した化学構造を外部のケモインフォマティクスツールに渡し、原子価や構造の妥当性を検証させ、その結果をフィードバックとして受け取りながらクレームを修正するという、決定論的で堅牢なシステムの構築が可能になる<sup>21</sup>。これは、ハルシネーションを本質的に排除するアプローチである。

## 5.2 無効審判とIRACフレームワークにおける高度な法的推論

特許の有効性や侵害を判断する高度な法的推論において、高性能LLMは目覚ましい適性を示す。「PILOT-Bench(Patent Invalidation Trial Benchmark)」は、米国特許審判部(PTAB)の査定系審判(Ex parte appeals)のデータセットを用い、LLMの特許領域における構造的な法的推論能力を評価するベンチマークである<sup>36</sup>。

PILOT-Benchは、法曹教育における標準的な法的推論フレームワークである「IRAC(Issue: 争点, Rule: 規範, Application: 当てはめ, Conclusion: 結論)」をモデル化している<sup>36</sup>。このベンチマークでは、PTABの審決文を分割し、審判請求人の主張(Appellant Arguments)と審査官の認定(Examiner Findings)という相反するテキストのみをモデルに入力する。モデルは、これらを統合的に分析し、以下のタスクを遂行しなければならない<sup>36</sup>。

1. **Issue Type**(争点の特定): 35 U.S.C.(米国特許法)の101条(特許適格性)、102条(新規性)、103条(進歩性)、112条(記載要件)などのうち、どの法的な拒絶理由が争点となっているかを特定するマルチラベル分類<sup>36</sup>。
2. **Board Authorities**(規範の特定): 争点に関連する37 C.F.R.(連邦規則集)などの手続き的・実体的な法的規範(ルール)を適用・マッピングする<sup>36</sup>。
3. **Subdecision**(結論の予測): 当事者の主張と法的規範をすり合わせた結果として、審判部の最終的な判断(Affirmed: 維持, Reversed: 取消, Affirmed-in-Part: 一部取消など)を予測するマルチクラス分類<sup>36</sup>。

低性能LLMが表面的なキーワード(例えば「103条」という単語の出現頻度)に引きずられて誤分類を起こすのに対し、高性能LLMは、「構成要件の認定」から「動機付けの有無」、そして「証拠の優越」に至る多段階の論理ステップ(Applicationのシミュレーション)を内部の思考連鎖(CoT)で展開できる<sup>36</sup>。これにより、弁理士がオフィスアクション(拒絶理由通知)に対応する際、審査官の論理の飛躍や法的な不備を指摘し、より強固で説得力のある反論(Response)を構築するための分析パートナーとして機能する<sup>37</sup>。

## 5.3 未知の技術領域(新規発明)へのゼロショット適応

特許の対象となる発明は、定義上「世界で初めて公開される新規な技術(Novelty)」である。したがって、過去の膨大な特許データセットをLLMに学習させても、真に革新的な発明の明細書を作成・評価するためには、過去のパターンの組み合わせ(結晶性知能)だけでは決して到達できない領域が存在する<sup>6</sup>。ここで求められるのは、未知の概念体系をその場で理解して論理を展開する「流動性

知能」である。

Gemini 3 Deep Think (ARC-AGI-2スコア 84.6%) やOpenAI o3 (同 75.7%~87.5%)、DeepSeek R1 といった最新モデルは<sup>14</sup>、事前の学習データに存在しない全く新しい視覚パズル (ARC-AGI-2) を解くと同様のプロセスで、新規発明に対応する。すなわち、プロンプトを通じて与えられた新規発明の開示書 (Invention Disclosure) や発明者との短い対話から、未知の技術的特徴とその因果関係 (なぜその構成が特定の課題を解決するのか) を「その場で」モデル化し、先行技術との差異 (進歩性) を精緻な言語で主張する能力である。これは、熟練した特許実務家が「発明者の頭の中にある新しい概念」をヒアリングし、特許法という独特の枠組み (ドメイン固有言語: DSL) に翻訳していく高度な認知作業をAIが模倣し始めていることを意味する<sup>20</sup>。強化学習によって内部的に特化言語 (DSL) を形成する能力を持つ推論モデルは、検証が厳密な法務・知財領域において、汎用モデルを遥かに凌駕する適応力を発揮する<sup>20</sup>。

## 6. 「知能のコスト (Cost of Intelligence)」と実務における戦略的最適配置

高性能LLMが知的財産業務において圧倒的な論理的優位性を持つことは疑いようがない。しかし、これらのモデルを現実の実務ワークフローに導入する際には、「知能のコスト (Cost of Intelligence)」という新たな経済的・インフラ的課題が浮上する。ARC-AGI-2のような難解な論理タスクを攻略するために用いられるテスト時計算 (推論時に多大な演算を行うアプローチ) は、計算コストと処理時間 (レイテンシ) を劇的に増大させるためである<sup>14</sup>。

### 6.1 A/E比 (精度対エネルギー比) と推論コストの非対称性

この「知能のコスト」を定量的に評価するための新たな指標として、金融分野の投資効率を示すP/E (株価収益率) の概念をAI分野に転用した「A/E (Accuracy/Energy: 精度対エネルギー比)」が提案されている<sup>39</sup>。MMLU (Massive Multitask Language Understanding) データセットを用いたクロスドメインのエネルギー消費ベンチマークの実証データによれば、LLMの推論にかかるエネルギーはタスクの複雑さによって大きく変動する<sup>39</sup>。医学や法務といった専門知識を要するドメインのクエリは、一般的な知識検索タスクと比較して、入力の長さや内部処理の複雑さから最大で4.3倍のエネルギー (およびそれに伴うAPIトークン消費コスト) を要求する<sup>39</sup>。

ARC-AGI-1の評価において、OpenAIのo3モデルが人間を上回る87.5%のスコアを達成した際、その「高計算モード (High Compute)」では探索とサンプリングに膨大なリソースを消費し、1タスクあたり推定3,400ドルという商業利用には非現実的なコストが発生したと試算されている<sup>14</sup>。より効率化された最新モデルであるGemini 3 Deep Think (ARC-AGI-2スコア84.6%) であっても、ARC Prize Foundationの検証によれば1タスクあたり13.62ドルの推論コストがかかっている<sup>14</sup>。一方で、DeepSeekのR1-Zeroは、強化学習による自己最適化の恩恵により、1タスクあたりわずか0.11ドルで14%のスコアを出し、コスト効率において異なる軌道を示している<sup>20</sup>。このように、論理的推論の深さとインファレンスコストは比例関係以上の指数関数的な増加を示すことがあり、すべての業務に一律に最高性能の推論モデルを適用することは、企業の知財部門や特許事務所にとって経済的破綻を

意味する<sup>41</sup>。

## 6.2 マルチホーミングとタスクの性質に応じたルーティング

この「知能のコスト」の非対称性は、知財実務において「マルチホーミング (Multihoming: 複数モデルの戦略的使い分けと同時利用)」の概念を必須のものとする<sup>43</sup>。業務の性質(要求される論理的深さとリスク許容度)に応じて、使用するAIモデルを動的にルーティングするアーキテクチャの設計が求められる<sup>41</sup>。

以下の表は、横軸に推論コスト(インファレンスタイムやAPI費用)、縦軸に要求される論理的深さ(ARC-AGI-2のスコアに相当する流動性知能)をとった場合の、知財タスクの最適な配置例を示している。

| 知財タスクの分類  | 要求される論理的推論レベル    | 許容されるインファレンスコスト | 最適なAIモデルの種類 | 具体的なユースケース  |
|-----------|------------------|-----------------|-------------|---|
| ルーチン型情報処理 | 低(表面的な言語処理、定型化)  | 低(高速・大量処理が必須)   | 低性能/ベースLLM  | 外国特許文献の翻訳、発明届出書の要約、定型的な書誌的事項の抽出、クライアント向けの一般的な進捗報告書の草案作成 |
| 中規模分析・検索  | 中(文脈理解と意味検索)     | 中(バッチ処理可能)      | 標準的LLM      | 先行技術文献の要約と技術分野の分類(IPC付与の補助)、類似特許のスクリーニング                |
| 高度な権利構築   | 高(シンボリック解釈、合成推論) | 高(精度がコストを正当化)   | 推論特化型LLM    | 特許請求の範囲(独立・従属クレーム)の論理構造構築、マーカッシュ構造の妥当性検証とサポート要          |

|           |                     |               |          |   |
|-----------|---------------------|---------------|----------|---|
|           |                     |               |          | 件の確認  |
| 批判的法的論理構築 | 最高(IRAC論理、コンテキスト適応) | 最高(致命的リスクの回避) | 推論特化型LLM | 拒絶理由通知(Office Action)に対する反論ロジックの組み立て、侵害予防調査(FTO)におけるクレーム解釈と対象製品の緻密な対比 |

システム設計の観点からは、分散型AIネイティブアーキテクチャやマルチエージェントシステムの導入が進行している。例えば、推論と全体計画を担う中核のオーケストレーター・エージェントには高コストな推論モデル(Reasoning LLMs)を配置し、特許データベースへのAPIコール、Markdown形式の文書整形、単純なデータバリデーションといった周辺タスク(ツール呼び出し)には、高速で安価な軽量モデル(例: GPT-4o-miniやオープンソースの小規模モデル)を割り当てるハイブリッド構成が、次世代の知財AIプラットフォームの標準構成となりつつある<sup>35</sup>。これにより、全体のパフォーマンスオーバーヘッドを抑制しつつ、致命的なエラー(Straggler Agentsによる失敗)を防ぐことが可能になる<sup>41</sup>。

### 6.3 弁理士法と倫理的・セキュリティ的考慮

AIを知財実務に統合する上で避けて通れないのが、法規制と倫理的要件への準拠である。低性能・高性能を問わず、パブリックなクラウドベースのLLMに発明の詳細を入力することは、そのデータがモデルの再学習に利用される可能性があり、特許法における「公知(Public Disclosure)」を構成し、新規性を喪失する重大なリスクをはらんでいる<sup>29</sup>。米国特許商標庁(USPTO)や欧州特許庁(EPO)、および日本弁理士会(JPAA)は、このリスクに対して強い警鐘を鳴らしており、エンタープライズグレードの機密保持契約や、データが学習に使用されないオプトアウト設定(あるいはローカル/プライベートクラウド環境での運用)を必須としている<sup>4</sup>。

さらに、日本においては弁理士法第75条(非弁理士の業務の禁止)との関係が議論されている<sup>4</sup>。AIを用いて特許明細書や特許庁への提出書類を完全に自動生成し、人間の弁理士の適切な監督や法的評価を介在させずにサービスを提供する事業者は、弁理士法に抵触する恐れがある。したがって、AIツールはあくまで「人間の実務家を補完し、思考を拡張するツール(Assistive Tools)」として位置づけられなければならない<sup>29</sup>。高性能LLMが提供する「説明可能性(Explainability)」や、生成された論理の出所を追跡可能な監査証跡(Audit Trails)機能は、実務家がAIの出力を法的に担保し、最終的な責任を負うための不可欠な要素となる<sup>21</sup>。

## 7. 結論

ARC-AGI-2ベンチマークが白日の下に晒した「言語生成能力と論理推論能力の決定的な乖離」は、

そのまま知的財産実務におけるAI利活用の分水嶺となる。単なる確率的なトークン予測に依存する旧来のベースモデル(低性能LLM)は、流暢なテキストを生成する能力には長けているものの、特許実務に不可欠な「深い流動性知能」――すなわち、厳密な法的シンボルの解釈、複雑なクレーム階層の合成的維持、および未知の技術コンテキストへの適応力――を欠いている。これらのモデルを、クレーム作成や侵害判断といったコア業務に単独で適用することは、先行詞の不一致や化学的ハルシネーションといった構造的欠陥を生み、特許の無効化という致命的なリスクを招く<sup>21</sup>。

一方で、テスト時計算と強化学習によって自己検証と論理探索の能力を獲得し、ARC-AGI-2において人間レベルに迫る飛躍的なスコアを達成した最新の推論特化型モデル群は、知財実務を「単なる文章の自動生成」から「論理制約の充足と自律的検証(Structure-First)」の次元へと引き上げつつある<sup>14</sup>。これらのモデルは、IRACフレームワークに基づく厳密な法的推論の展開や、未知の新規発明の概念モデル化において、人間の弁理士に匹敵する、あるいはそれを高度に補完する知的パートナーとなり得る<sup>36</sup>。

しかしながら、この高度な推論能力には、計算リソースとレイテンシという明確な「知能のコスト」が伴う<sup>39</sup>。今後の知的財産テクノロジー(IP Tech)における競争優位性は、単一のAIモデルにすべてを委ねる技術的盲信ではなく、タスクの性質とリスク許容度に応じて、ルーチンワークを担う「低コストな汎用LLM」と、致命的リスクを検証する「高コストな推論LLM」を動的にオーケストレーションするアーキテクチャ設計の巧拙によって決定づけられる。

知的財産実務家は、AIを単なる「執筆の代替(ライティングツール)」としてではなく、「論理構築と検証の協働者(リーズニングパートナー)」として再定義する必要がある。人間の専門家が有する法的な倫理観や高度な戦略的判断力と、高性能LLMの流動性知能を適切に融合させることでのみ、縮小するリソース環境下においても特許の権利価値を最大化し、グローバルなイノベーション競争を勝ち抜くことが可能となる。

## 引用文献

1. 弁理士業界の将来とAIの利用, 2月 22, 2026にアクセス、  
<https://jpaa-patent.info/patent/viewPdf/4642>
2. 各小委員会の報告 - 特許庁, 2月 22, 2026にアクセス、  
[https://www.jpo.go.jp/resources/shingikai/sangyo-kouzou/shousai/chizai\\_bunkakai/document/19-shiryuu/04.pdf](https://www.jpo.go.jp/resources/shingikai/sangyo-kouzou/shousai/chizai_bunkakai/document/19-shiryuu/04.pdf)
3. 日本弁理士会協賛セッション, 2月 22, 2026にアクセス、  
<https://www.ipaj.org/workshop/2025/pdfs/S03.pdf>
4. AI等を用いた業務支援サービスの提供と弁理士法第75条との関係, 2月 22, 2026にアクセス、  
<https://www.jpaa.or.jp/cms/wp-content/uploads/2025/04/AIservices-article75.pdf>
5. ARC-AGI: Benchmark for Abstraction & Reasoning in AGI, 2月 22, 2026にアクセス、  
<https://www.emergentmind.com/topics/abstraction-and-reasoning-corpus-arc-agi>
6. ARC-AGI-2: A New Challenge for Frontier AI Reasoning Systems, 2月 22, 2026にアクセス、  
<https://arxiv.org/html/2505.11831v2>
7. Can LLMs Reason?. Attempting the ARC-AGI challenge with, 2月 22, 2026にアクセス

- ス、  
<https://medium.com/bitgrit-data-science-publication/can-llms-reason-c58a45e17059>
8. Is Your AI Smart Enough? Test It with ARC AGI v2! - Labellerr, 2月 22, 2026にアクセス、<https://www.labellerr.com/blog/arc-agi-v2/>
  9. The Evolving Landscape of AI Intelligence: A Comparative Analysis, 2月 22, 2026にアクセス、  
<https://medium.com/ai-simplified-in-plain-english/the-evolving-landscape-of-ai-intelligence-a-comparative-analysis-of-arc-agi-1-and-arc-agi-2-9e8782222c8d>
  10. (PDF) Problem-Solving and Intelligence - ResearchGate, 2月 22, 2026にアクセス、  
[https://www.researchgate.net/publication/368088357\\_Problem-Solving\\_and\\_Intelligence](https://www.researchgate.net/publication/368088357_Problem-Solving_and_Intelligence)
  11. US8821242B2 - Systems and methods for enhancing cognition, 2月 22, 2026にアクセス、<https://patents.google.com/patent/US8821242B2/en>
  12. GPT-5.2 & ARC-AGI-2: A Benchmark Analysis of AI Reasoning, 2月 22, 2026にアクセス、<https://intuitionlabs.ai/articles/gpt-5-2-arc-agi-2-benchmark>
  13. ARC-AGI-2 - ARC Prize, 2月 22, 2026にアクセス、<https://arcprize.org/arc-agi/2/>
  14. Announcing ARC-AGI-2 and ARC Prize 2025, 2月 22, 2026にアクセス、  
<https://arcprize.org/blog/announcing-arc-agi-2-and-arc-prize-2025>
  15. Product of Experts with LLMs., 2月 22, 2026にアクセス、  
<https://lambdalabsml.github.io/Product-of-Experts-ARC-Paper/>
  16. Guide - ARC Prize, 2月 22, 2026にアクセス、<https://arcprize.org/guide>
  17. We tested every major AI reasoning system. There is no clear winner., 2月 22, 2026にアクセス、<https://arcprize.org/blog/which-ai-reasoning-model-is-best>
  18. Gemini 3 Deep Think Achieves 45.1% on ARC-AGI-2 - Reddit, 2月 22, 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/)
  19. Google Gemini 3 Deep Think Hits 84.6% on ARC-AGI-2, Beating, 2月 22, 2026にアクセス、  
<https://www.implicator.ai/google-gemini-3-deep-think-hits-84-6-on-arc-agi-2-beating-gpt-5-and-claude-2/>
  20. An Analysis of DeepSeek's R1-Zero and R1 - ARC Prize, 2月 22, 2026にアクセス、  
<https://arcprize.org/blog/r1-zero-r1-results-analysis>
  21. AI Patent Drafting Tools in 2026: Can They Handle Chemistry, 2月 22, 2026にアクセス、  
<https://www.deepip.ai/blog/ai-patent-drafting-tools-chemistry-biotech-2026>
  22. Multi-dimensional Evaluation of LLM-Generated Patent Claims - arXiv, 2月 22, 2026にアクセス、  
<https://arxiv.org/html/2505.19345v2>
  23. The Role of LLMs/ Generative AI in Patent Litigation - XLSCOUT, 2月 22, 2026にアクセス、  
<https://xlscout.ai/the-role-of-large-language-models-in-patent-litigation/>
  24. Multi-dimensional Evaluation of LLM-Generated Patent Claims, 2月 22, 2026にアクセス、  
[https://www.researchgate.net/publication/392104195\\_PatentScore\\_Multi-dimensional\\_Evaluation\\_of\\_LLM-Generated\\_Patent\\_Claims](https://www.researchgate.net/publication/392104195_PatentScore_Multi-dimensional_Evaluation_of_LLM-Generated_Patent_Claims)
  25. ARC 2 looks identical to ARC 1. Humans get 100% on it (early results)., 2月 22,

- 2026にアクセス、  
[https://www.reddit.com/r/singularity/comments/1j1ao3n/arc\\_2\\_looks\\_identical\\_to\\_arc\\_1\\_humans\\_get\\_100\\_on/](https://www.reddit.com/r/singularity/comments/1j1ao3n/arc_2_looks_identical_to_arc_1_humans_get_100_on/)
26. IPBench: Benchmarking the Knowledge of Large Language Models, 2月 22, 2026  
にアクセス、<https://arxiv.org/html/2504.15524v2>
  27. Unraveling the Perils and Promises in Patent Drafting - TT Consultants, 2月 22,  
2026にアクセス、  
<https://ttconsultants.com/chatgpt-as-a-double-edged-sword-unraveling-the-perils-and-promises-in-patent-drafting/>
  28. Transforming Patent Litigation with AI: The Strategic Role of LLMs, 2月 22, 2026に  
アクセス、  
<https://xlscout.ai/transforming-patent-litigation-with-ai-the-strategic-role-of-llms-and-generative-ai/>
  29. The hidden risks of using LLMs in the inventive process - GJE, 2月 22, 2026にアク  
セス、  
<https://www.gje.com/resources/the-hidden-risks-of-using-llms-in-the-inventive-process/>
  30. Towards Automated Quality Assurance of Patent Specifications - arXiv, 2月 22,  
2026にアクセス、<https://arxiv.org/pdf/2510.25402>
  31. IPBENCH - OpenReview, 2月 22, 2026にアクセス、  
<https://openreview.net/pdf/aeb1d653da4f45e1403b114cccfd66efb69e6c56.pdf>
  32. IPBench, 2月 22, 2026にアクセス、<https://ipbench.wangqiyao.me/>
  33. IPBench: Benchmarking the Knowledge of Large Language Models, 2月 22, 2026  
にアクセス、  
[https://www.researchgate.net/publication/391020110\\_IPBench\\_Benchmarking\\_the\\_Knowledge\\_of\\_Large\\_Language\\_Models\\_in\\_Intellectual\\_Property](https://www.researchgate.net/publication/391020110_IPBench_Benchmarking_the_Knowledge_of_Large_Language_Models_in_Intellectual_Property)
  34. 「集合知」と「試行錯誤」によるフロンティアAIの推論時スケーリング, 2月 22, 2026にアク  
セス、<https://sakana.ai/ab-mcts-jp/>
  35. SotA ARC-AGI-2 Results with REPL Agents | Symbolica Blog, 2月 22, 2026にアクセ  
ス、<https://www.symbolica.ai/blog/arcgentica>
  36. PILOT-Bench: A Benchmark for Legal Reasoning in ... - ACL Anthology, 2月 22,  
2026にアクセス、<https://aclanthology.org/2025.nllp-1.17.pdf>
  37. Best AI Tools for Patent Attorneys 2026: Drafting, Search & Office, 2月 22, 2026に  
アクセス、  
<https://www.deepip.ai/blog/top-ai-tools-patent-attorneys-are-using-to-boost-efficiency>
  38. AI Tools for Patent Drafting: LLMs Will Likely Never Write Claims as, 2月 22, 2026に  
アクセス、  
<https://ipwatchdog.com/2024/04/21/ai-tools-patent-drafting-llms-will-likely-never-write-claims-well-humans/>
  39. The Carbon Cost of Intelligence: A Domain-Specific Framework for, 2月 22, 2026  
にアクセス、<https://www.mdpi.com/1996-1073/19/3/642>
  40. Chapter 14: The World's 1st On-Chain LLM | by Modulus Labs, 2月 22, 2026にアク  
セス、  
<https://medium.com/@CountableMagic/chapter-14-the-worlds-1st-on-chain-llm->

[7e389189f85e](#)

41. An Open-Source White-Box Agentic Benchmark Suite for AI-Native, 2月 22, 2026  
にアクセス、<https://arxiv.org/html/2601.09393v1>
42. The Probability Landscape, 2月 22, 2026にアクセス、  
<https://papers.ssrn.com/sol3/Delivery.cfm/6001374.pdf?abstractid=6001374&mirid=1>
43. Stop using the wrong AI platforms: LLMs vs. Reasoning Models, 2月 22, 2026にアクセス、  
<https://punchcard.io/blog/stop-using-the-wrong-ai-platform-llms-vs-reasoning-models/>
44. Pricing, Supply, and Demand for LLMs - Andrey Fradkin, 2月 22, 2026にアクセス、  
[https://andreyfradkin.com/assets/LLM\\_Demand\\_12\\_12\\_2025.pdf](https://andreyfradkin.com/assets/LLM_Demand_12_12_2025.pdf)
45. Reasoning, tool calling, and agentic systems in legal LLMs, 2月 22, 2026にアクセス、  
<https://legal.thomsonreuters.com/blog/reasoning-tool-calling-and-agentic-systems-in-legal-llms/>
46. Can Large Language Models Understand As Well As Apply Patent, 2月 22, 2026にアクセス、  
<https://arxiv.org/html/2507.10576v1>