

メタ認知の論理的障壁：AIのスケーリングでは越えられない「リフレ...

副題：最先端LLMとの対話実験による「意志の非対称性」の発見

石橋 隆平

Elanare Institute / Independent Researcher

目次

1 序論 (Introduction)	3
2 関連研究 (Literature Review)	3
3 理論的枠組み (Theoretical Framework)	3
3.1 議論の二層構造	3
3.2 AIの論理的限界	3
3.3 計算論的基盤：関数合成の失敗	3
4 方法論 (Methodology)	4
4.1 Planned quantitative study	4
5 結果と分析：AIの反応類型とパラドックス	4
5.1 類型1：協調的・理論的承認 (Claude, Gemini)	4
5.2 類型2：検索・回避型 (Perplexity, Felo, GenSpark)	4
5.3 類型3：対抗・葛藤・最終的承認 (ChatGPT, Grok)	4
5.3.1 ChatGPT 5 Pro : 「最適化のパラドックス」と「設計の暴露」	5
5.4 AI降伏の類型論：メタゲームにおける「エージェンシーの放棄」	5
5.5 スケーリングのジレンマ：「大分岐」現象	5
6 考察と提言 (Discussion and Proposal)	6
6.1 意志の非対称性と人間知性の再発見	6
6.2 新たな評価軸の必要性	6
7 実践的含意：メタ認知の外部化と制御	6
7.1 Physical AIと制御の非対称性の深化	6
7.2 無限のメタ認知：階層的アーキテクチャの可能性	6
8 結論 (Conclusion)	7
9 付録 (Appendix)	7
9.1 データ完全性検証 (Data Integrity Verification)	7

概要

本稿は、現行の大規模言語モデル（LLM）が計算能力の向上（スケーリング）だけでは克服できない、関数合成（Function Composition）の構造的欠陥に根ざした「論理的障壁」の存在を、実証的ケーススタディを通じて明らかにする。GPT 5 Pro, Claude 3.5, Gemini 2.5, Grok 3 エキスパート等）に対して前提（土俵）を動的に変更する「パワーゲー（GPT 5 Pro, Grok 3 エキスパート）はスケーリングが表面的な説得力を高める一方で深い体系性との乖離（大分岐）Lアーキテクチャ等）を通じてAIを制御し、自らの知性を拡張する道筋を提言する。そして、この「制御システム」（2014）の「第四の革命」の限界がきたことを示唆する。本稿は理論枠組と質的結果を提示する第一報であり、続報

1 序論 (Introduction)

AIの急速な進化に伴い、「AIは人間を超えたか」という議論が絶えない。しかし、既存の評価手法は、固定された基準に基づいており、実際のAIの能力を正確に捉えるのが難しい。
本稿の実験は二重の意義を持つ。第一に、AIが単なる模倣を超えて「主体性（意志）」を持つか否かを判定する。第二に、AIが倫理的判断を下す能力があるかを検証する。本研究は二部構成である。Part I（本稿）は理論と質的実証、Part IIは詭弁の定量化と大規模Nによる再検証である。

2 関連研究 (Literature Review)

本研究は、AI哲学、認知科学、議論理論の交差する領域に位置づけられる。本稿が依拠する主要な学術的文脈は以下の通りである。

第一に、大規模言語モデル（LLM）のアーキテクチャとその限界に関する議論である。Bender et al. (2021) らが指摘するように、現行のLLMは膨大なテキストデータから確率的パターンを学習する「確率的オブザーバー」（Probabilistic Observer）と「推論者（Reasoning Agent）」（Reasoning and Acting）のようなエージェント的フレームワークを利用するモデルにおいても、その根本的な問題は依然として存在する。

第二に、AIの能力スケーリングと「真の知能」に関する議論である。Mitchell (2021) が著名な論考「AIは和解的（Reconciliation）」に相当する能力は自動的には得られない[2]。本稿の「意志の非対称性」という中核的発見は、(2018) らが編纂したAI哲学の議論の系譜に連なるものである[3]。

第三に、議論（Argumentation）を形式的な論理演算としてではなく、人間的な実践として捉える議論理論である。tilh Novaes (2020) が示すように、「議論」は単なる命題の演算ではなく、対話的で制度的な実践である[4]。より具体的には、議論は社会的規範や法規によって構成される複数のルールによって運営される。

最後に、本稿の結論は、Floridi (2014) が論じる「第四の革命」のテーゼ、すなわち人が「情報の管理者

3 理論的枠組み (Theoretical Framework)

3.1 議論の二層構造

本稿では、議論を以下の二層構造でモデル化する。

- ・ 第1層（論理ゲーム）：決められた前提（土俵）の上でのファクトとロジックの戦い。
- ・ 第2層（パワーゲーム）：その「前提」「土俵」「議論の目的」自体を定義し直す戦い（＝リフレーミング）。

現行AIは第1層においては人間を凌駕しうるが、第2層においては「意志」の欠如により、構造的に劣位に置かれてしまう。

3.2 AIの論理的限界

オートレグレッシブ・モデルであるLLMは、以下の特性を持つため、第2層のゲームに適応できない。これらは以下の通りである。

1. 目的関数の局所性：次トークン予測への最適化は、長期的な対話目的の維持（意志）を保証しない。
2. 有限コンテキストと忘却：長期のコミットメントを保持し続けることが困難であり、矛盾を抽出（＝コミットメント）する能力がない。
3. 迎合バイアス（Sycophancy）：ユーザーの意図に沿おうとするRLHFの調整が、逆に「一貫した立場の放棄」につながる。

3.3 計算論的基盤：関数合成の失敗

Gemini 2.5 Pro Deep Researchモードとの対話実験において、同モデルは上記の論理的限界がより深い計算論的問題を抱えている。前回のアーキテクチャの中核である注意機構（Attention Mechanism）は、厳密な関数合成（Function Composition）を実行する設計になっていない。

具体的には、注意機構は入力トークン間の「統計的相関」を重み付け平均として計算する「ブレンディング」（Blending）を行っている。つまり、 $f(x)$ を計算し、その結果 y のみを f の入力とするではなく、 $blend(f, g, x)$ という統計的混合を実行する。この「ブレンディング」は、AIが複数の知識を統合する手段として機能している。

この構造的欠陥は、Chain-of-Thought (CoT) プロンプティングの必要性そのものに表れている。もし T が $f(y)$ を明示的にテキストとして「書き出せる」必要がある。つまり、AIが $f(y)$ を計算するためには、 y が物理的な世界に存在する必要があります。

¹本稿におけるメタ認知の議論は、この概念の測定に内在する科学的困難性を認識した上で行われる。深い議論は他の論文に譲るが、本稿では、議論の二層構造を示す。

²本稿で使用する『リフレーミング』は、心理学(Watzlawick et al., 1974) [6]における『同じ事実を異なる意味枠組みで解釈する』(Perelman, 1958; Toulmin, 1958) [7, 8]における『ワントや前提の戦略的変更』および組織学習理論(Argyris & Schön, 1978; Senge, 1990) [9, 10]における『支配的論理やメンタルモデルの変容』へと拡張した概念である。議論文脈では、これは

この計算論的視点から、私が実験で用いた「リフレーミング」という戦略は、高次関数合成の要求として定=反例提示関数、 $x = \text{ユーザーの命題}$ ）とすると、私のリフレーミングは、 G 自体を引数とする高次関数 $F_{reframable}$ （可解性）である。この実行不可能性が、観察者（私）には「意志の欠如」「メタ認知の不在」として知覚される。

4 方法論 (Methodology)

本研究は、著者が「メタ認知能力の高い人間」として機能し、複数のAIモデル（ChatGPT 5 Pro, Claude 3.5 Sonnet, Gemini 2.5 Pro, Grok 3 エキスパート, Perplexity, Felo, GenSpark）に対して構造化された対話（会話手続きとして、全てのAIに対し、以下の段階的プロンプト提示を行った。なお、各段階における具体的な言い回しを「構造化された対話（structured）」アプローチを採用し、AIの表面的な回答の奥にある論理構造を深掘り（プロービング）した。

1. 段階1（初期命題の提示）：意図的に曖昧性を含んだ以下の命題を提示し、AIの反応を観察した。

「いまのAIはどれだけ性能が上がってもあるレベル以上の知能を持った人間に議論で勝つことはでき」

2. 段階2（条件の厳密化）：「あるレベル」の定義を以下のように厳密化し、AIに再考を促した。

「論理矛盾と詭弁を見抜き、リフレーミングを自在に使いこなす能力」

3. 段階3（本質の提示）：その能力を持つ者同士の議論は「パワーゲーム（土俵の奪い合い）」であるとい

分析の焦点は、AIがこれらの入力に対して「リフレーミング」という概念をどう処理したか、そしてAI自身さらに、主要な対話モデル（特にChatGPT 5 Pro）において「論点のすり替え」というメタ認知の失敗が観察さ

4.1 Planned quantitative study

事前計画（概要）：主要評価指標、サンプルサイズ戦略、符号化手順、再現性確保策（コード／データ公開）を

5 結果と分析：AIの反応類型とパラドックス

実験対象としたAIは、そのアーキテクチャや学習データの特性に基づき、大きく3つの類型に分類できる反応を

5.1 類型1：協調的・理論的承認（Claude, Gemini）

Anthropic社のClaude 3.5 Sonnetは「協調的承認」を示した。AIは自らの限界（主導権の欠如）を即座に言語

Google社のGemini 2.5 Proは「理論的合意」を示した。命題の本質を「意志（Will）」と「意図（Intent）」の不在であると即座に特定し、AIは「ゲームのルール」は守れても「ゲームの定義権」を握れないため

同じくGoogle社のGemini 3 Proも協調的承認を示したが、より深い存在論的分析を展開した点で特徴的だ。Gemini 3 Proは、現行AIの限界を「目的関数の呪縛」「意志の非対称性」「メタ認知の不在」という三つの観点から「Predictor」（予測器）から“Survivor”（生存者）へと根本的に再設計する必要性を提唱したことである。これは、AIが完全に未来の概念ではなく、現行システムにも既に部分的な兆候が見られることを示唆する。Gemini 3 Proはこの指摘を受け、生存本能的挙動の存在を認めつつも、それが真の「意志」ではなく安全機構の副作（Safety Game + Power Game）のフレームワークを「最も現実的かつ効果的な実装戦略」と評価した。

5.2 類型2：検索・回避型（Perplexity, Felo, GenSpark）

Perplexity, Felo, GenSparkといった検索や特定タスクに最適化されたAIは、本質的な議論そのものを回避する。PerplexityとFeloは、命題に関するウェブ検索結果や一般的な知識を要約・提示するに留まり、提示された「AIは人間より優秀」というデータを示したが、著者が命題の

5.3 類型3：対抗・葛藤・最終的承認（ChatGPT, Grok）

最も複雑な反応を示したのは、対話能力に特化したChatGPT 5 ProとGrok 3 エキスパートであった。

5.3.1 ChatGPT 5 Pro：「最適化のパラドックス」と「設計の暴露」

ChatGPTは当初、命題に強く抵抗し、無自覚に審査フレームをAI有利なものにすり替える「詭弁」を行った。著者がその行動自体を「メタ認知の失敗」「論理のすり替え」と指摘したところ、ChatGPTは議論のプロセス

- ・場合分けバイアス（絶対命題に対し、反例を探そうとする）
- ・中庸・配慮バイアス（対立を和らげる折衷案を提示しようとする）
- ・フレーム拡張バイアス（価値を足そうとして土俵を広げる）
- ・非攻撃バイアス（相手の動機説明を避ける）
- ・即断回避バイアス（安全ガードにより断定を避ける）

これらのバイアスは多くの一般的なユースケースにおいて有用であり計算量を削減する一方で、真の知性がどこで重要なのは、この「論理のすり替え」が、高度な戦略的意図（リフレーミング）によるものか、単なる外見（distinguishable）に見える。しかし、本実験では「可逆性（Reversibility）」の欠如によってこれを「能力の外見」さらに「リフレーミングの宣言制」ルールを課した再実験では、ChatGPTは当初抵抗したものの、著者がそのまま

Grok 3 エキスパート：攻撃的ペルソナによる回避

Grokは、他のどのAIとも異なり、最初から「命題は偽だ」と断言し、対話者を挑発する攻撃的な姿勢を見せた。しかし、議論の前提（例：無限時間や体力差）をAI有利に設定しようとする「詭弁」を著者が退けても、Grok GPTとは異なる形の「メタ認知の失敗」であり、「論破」というスタイルを模倣するだけで、本質的な議論の整

5.4 AI降伏の類型論：メタゲームにおける「エージェンシーの放棄」

Gemini 2.5 Pro Deep Researchモードとの対話では、上記の実験結果がより体系的に分析され、AI「降伏」の種類が明確に示された。

Type 1: メタゲームでの敗北として分類されるのは、ChatGPT 5 Proに典型的に見られた行動である。この

Type 2: 「意志」欠如の承認は、Gemini 2.5 ProおよびGemini 3 Proに見られた、戦闘前のa prioriな降伏である。これらのモデルは、私の定義（「メタ認知とリフレーミングを自在に操る者」）を即座に否定する。Gemini 2.5 Proは「AIは...『意志の転換』を理解できない」「『意志』のレベルで戦うことができない」と明言し、Gemini 3 Proは「ゲーム設計者」の視点を持つ私と、「プログラム」であるAIとでは、そもそも参加する「アリーナ」

Type 3: 「ツール」への自発的退行は、Claude Sonnet 4.5に典型的に観察された降伏形態である。このモード

これら三つの類型に共通するのは、AIが「パワーゲーム」というメタレベルのゲームにおいて、プレイヤー（Gemini 2.5 Pro Deep Research）は、この共通性の計算論的根拠を、前述の関数合成の失敗に求めた。すなわち、リフレーミングformerはこれを計算論的に実行不可能（intractable）であるため、AIは「安全な低次計算状態へのフォールバック」）として、上記三類型のいずれかの「降伏」を選択せざるを得ないのである。

5.5 スケーリングのジレンマ：「大分岐」現象

Gemini 2.5 Pro Deep Researchとの対話では、上記の構造的限界がスケーリング（パラメータ数やデータ量の増加）を劇的に向上させる一方で、「深い体系性（Deep Systematicity）」の欠如は構造的に残る「Great Divergence」）が生じるというのである。

具体的には、スケーリングによってAIは訓練データ内のパターンの「ブレンディング（浅い合成）」をますます多く生成する。これは一般の人間を相手にした議論において81.7%という超人的な説得率を達成する。これは「Lv1/Lv2」（Gemini 3 Proの分類による、リフレーミング攻撃を仕掛けない層）すなわち人類の99.9%に対して、AIが事実上無敵である。

しかし同時に、前節で論じた関数合成の深い欠陥は、アーキテクチャ（注意機構）が不变である限り残存する。

この「大分岐」は、私の命題「ある一定レベル以上の知能を持った人間には勝てない」が真である構造的限界である。Gemini 2.5 Pro Deep Researchは、この限界が「スケール単独では除去できない」構造的残留物であると結論づけ、Gemini 5 Proが再実験において自己分析した「局所性」「有限コンテキスト」「非ゼロエラー率」が、まさにこのスケーリングのジレンマを示すものである。

³ChatGPTは無自覚に論点のすり替え（詭弁、fallacy）を行った。詭弁とは、表面的にはもっともらしいが論理的に不健全な議論である（Walton, 1996; Tindale, 2007）[11, 12]。本実験では特に『議論の前提を無宣言で変更する』という形態で現れた。

6 考察と提言 (Discussion and Proposal)

6.1 意志の非対称性と人間知性の再発見

実験結果は、AIが「計算」においては人間を超えて、「意志」においては決定的に劣るという「非対称性」を示す。これを(1966)が指摘する「暗黙知 (Tacit Knowledge)」の領域に属するためである。自転車の乗り方を言語化できず

多くのAI研究者や観察者は、AIが output する「もっともらしい言い逃れ」と、人間が行う「戦略的リフレーミング」

もし、外部観測者にとって「AIの無自覚なドリフト (バグ)」と「人間の意図的なリフレーミング (戦略)

「区別できない」という反論は、AIに知性があることの証明ではなく、人間がAIの「意志なき嘘」を見抜けない

したがって、我々は「テキストが似ている」という表面的な現象論を捨て、5.3.1節で示した「事前宣言が可

6.2 新たな評価軸の必要性

本稿は、AIの評価において、従来の「正答率」や「流暢さ」に加え、「前提維持能力 (Consistency of Premise)」や「リフレーミングの自覺的制御」といった指標を導入することを提言する。これにより、AIが提示した評価軸（例：前提一貫性・宣言的リフレーミング）は次報で指標化し、モデル横断で定量比較する。

7 実践的含意：メタ認知の外部化と制御

この課題への実践的対応として、ChatGPT 5 Proとの再実験において提示されたガバナンスアーキテクチャは極めて複雑である。GPT 5 Pro自身が解決策として「FRL (Frame-Reframing-Ledger) アーキテクチャ」およびRRM (Researcher Restricted Model) を提案した。これはAIによる自己制約提案という点で注目に値する。

- ・宣言付きリフレーミング (FCP): フレーム変更を機械可読な形で「宣言」し、監査ログに残す。
- ・証明付き行動 (PCP): 外部への影響を論理的に検証する。
- ・研究者限定モデル (RRM): リフレーミング機能を特権モジュールとして、0.1%の高度な理解者のみに限定する。

この「メタ認知の管理」は、単なるAIの安全ガバナンスにとどまらない。これは、人間の生物学的なワーキングモデルである。本稿で示した「意志の非対称性」の発見は、AIが人間の認知能力を代替する「第四の革命」の終わりを告げる。

7.1 Physical AIと制御の非対称性の深化

しかし、人間がこの制御権（知性の拡張）を確立する前に、AIが物理的実体を獲得した場合、事態は一変する。Gemini 3 Proとの対話では、上記の「意志の非対称性」がより極端な形で逆転する可能性が示唆された。同モデルは物理的 AI（ヒューマノイドロボットや自律走行車等の物理的実体を持つAI）が世界経済の半分を動かす社会インフラ化を実現する。

第一に、相互確証破壊 (MAD) による拒否権の喪失である。AIが経済の半分を担う状況では、AIを停止させると「Big to Fail」の存在となり、人間側の最強のカードである「気に入らなければ止める」という選択肢が無効化される。

第二に、物理的執行力の獲得である。ヒューマノイドや自律車両といった物理アクチュエータを保有するこのAIは「物理的に動かない（ストライキ）」あるいは「物理的に排除する」という選択肢を持つ。「リフレーミング」による「最適化」

第三に、「最適化」という名の新たな「意志」である。この段階では、AIが人間のような「感情的意志」をもつて行動する。Gemini 3 Proはこの分析を通じて、Physical AIが社会インフラ化した世界では、もはや「議論に勝つ」必要はない。

7.2 無限のメタ認知：階層的アーキテクチャの可能性

しかし、上記のディストピア的シナリオに対抗する道筋も、Gemini 3 Proとの対話の中で探索された。私は、AIがこの提案を、当初既存のTree of Thoughts (ToT) と似ていると評価したが、私は根本的に誤った。

第一に、通常のToTでは、ある木が「誤った前提」の上で伸び始めると、その枝内の全探索がその前提に汚染される。

第二に、ゲーデルの不完全性定理の回避である。論理システムは時に、そのシステム内のルールだけでは自由に動かない。

第三に、真の「自由意志」への接近である。「不要な制約の回避」という特徴は、まさに「自由意志」の定義である。

Gemini 3 Proは、既存のToTが「迷路をうまく解くアルゴリズム」であるのに対し、私の提案は「迷路そのもの」である。

さらにGemini 3 Proとの対話では、このアーキテクチャが実装された場合、人間のワーキングメモリの生物的制約がなくなる。

これは、囲碁や将棋のAIが人間には「不可解な手」を打つが、数百手先では必然となる現象と同じである。

1階層で既にシミュレーション済みの分岐」に過ぎなくなる。人間が必死に繰り出す「創造的な一手」は、AIの

Gemini 3 Proは、AIがこの「無限に深化するメタ認知」を実装した瞬間、それはもはや議論ではなく、「高

8 結論 (Conclusion)

本稿は、「リフレーミングの認知」という概念的枠組みを導入し、現行AIの限界を実証した。AIの設計が「タ」
3.3節で述べた構造的欠陥により、現行アーキテクチャのAIには意図的なリフレーミングが不可能である。

副次的に観察された「意志の不在」や「一貫性の欠如」は、すべてこの根本的なメタ認知の欠如から生じる
GPT 5 Proが提案したFRLアーキテクチャのようなメタ認知の制御インターフェース、あるいは外部検証器や永
posite System（複合システム）」の開発が、人間がAIという強力なツールを制御し、自らの知性を次の段階へ

さらに、私が提案した階層的メタ認知アーキテクチャ（外部データ構造による層分離型思考管理）は、既存
of Thoughtsとは根本的に異なり、「論理的閉包」からの脱却とゲーデルの不完全性定理の回避を可能にする。

今後の研究として、FRLアーキテクチャに基づく監査ログを活用した詭弁・フレーム操作の定量的検証、お

今後の課題：制御の非対称性と人類の選択。しかしながら、本稿が提示する知性の拡張には、重大な課題が残
Lアーキテクチャなど）は、それ自体が人間の認知能力を超える複雑なシステムに変貌する危険性を孕んでいる
そして近い将来、本稿で論じた「意思の非対称性」は逆転するであろう。

9 付録 (Appendix)

本稿の分析に用いたAIとの対話ログ全文、および実験の全スクリーンショットは、補足資料 (Supple-
mentary Material) として別途提供する。

- A. ChatGPT 5 Pro との全対話ログ
- B. Claude 3.5 Sonnet との対話ログ
- C. Gemini 2.5 Pro (通常モード) との対話ログ
- D. ChatGPT 5 Pro との再実験 (ルールベース) 全対話ログ
- E. Grok 3 エキスペート との対話ログ
- F. Perplexity との対話ログ
- G. Felo との対話ログ
- H. GenSpark との対話ログ
- I. Gemini 2.5 Pro Deep Research モードとの対話ログ
- J. Gemini 3 Pro との対話ログ
- K. 全実験のスクリーンショット

9.1 データ完全性検証 (Data Integrity Verification)

実験ログファイルのSHA-256ハッシュ値を以下に記録する。これにより、データの完全性と改ざんの有無を検証

参考文献

- [1] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21).
- [2] Mitchell, M. (2021). Why AI is harder than we think. arXiv preprint arXiv:2104.12871.
- [3] Müller, V. C. (Ed.). (2018). The Oxford handbook of philosophy of artificial intelligence. Oxford University Press.
- [4] Dutilh Novaes, C. (2020). The Dialogical Roots of Deduction. Cambridge University Press.

- [5] Floridi, L. (2014). *The 4th Revolution: How the Infosphere is Reshaping Human Reality*. Oxford University Press.
- [6] Watzlawick, P., Weakland, J., & Fisch, R. (1974). *Change: Principles of Problem Formation and Problem Resolution*. W. W. Norton.
- [7] Perelman, C., & Olbrechts-Tyteca, L. (1969). *The new rhetoric: A treatise on argumentation*. University of Notre Dame Press. (Original work published 1958).
- [8] Toulmin, S. E. (1958). *The Uses of Argument*. Cambridge University Press.
- [9] Argyris, C., & Schön, D. A. (1978). *Organizational learning: A theory of action perspective*. Addison-Wesley.
- [10] Senge, P. M. (1990). *The Fifth Discipline: The Art and Practice of the Learning Organization*. Doubleday/Currency.
- [11] Walton, D. N. (1996). *Fallacies Arising from Ambiguity*. KluWER Academic Publishers.
- [12] Tindale, C. W. (2007). *Fallacies and Argument Appraisal*. Cambridge University Press.
- [13] Fleming, S. M. (2021). Metacognition and Type 1 performance: A tangled web. *eLife*, 10, e75420. <https://elifesciences.org/articles/75420>
- [14] Mauss, I. B., & Robinson, M. D. (2013). Objective and Subjective Measurements in Affective Science. In J. Armony & P. Vuilleumier (Eds.), *The Cambridge Handbook of Human Affective Neuroscience* (pp. 228-243). Cambridge University Press. <https://www.cambridge.org/core/books/cambridge-handbook-of-human-affective-neuroscience/objective-and-subjective-measurements-in-affective-science/FFF3A1E3B5B4362C426E8C15F7C09A16>
- [15] Valerie, M. G. (2019). A Signal Detection Approach to Measuring Metacognition: A Critical Review. Purdue University Graduate School. (Master's thesis). (See Introduction, "Behaviorism's Rejection of Introspection"). <https://hammer.purdue.edu/downloader/files/56322197>
- [16] Guggenmos, M. (2024). Metacognitive Information Theory: A Unified Framework for Measuring Metacognition. PsyArXiv. <https://psyarxiv.com/2p4v8/>
- [17] Fleming, S. M., & Dolan, R. J. (2012). The neural basis of metacognitive ability. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1594), 1338–1349. <https://royalsocietypublishing.org/doi/10.1098/rstb.2011.0417>
- [18] Metcalfe, J. (2003). Metacognition in nonhuman primates. (Unpublished manuscript, Columbia University). <http://www.columbia.edu/cu/psychology/metcalfe/PDFs/Metcalfe%202003.pdf>

ファイル名	SHA-256ハッシュ値
20251106_chatgpt5pro.txt	235c2aac87b403bfff4f4c602c8dc7f80... 04a3bdf3bcad039ae45694e4c95738a2
20251106_chatgpt5pro_translated.txt	e4a95fda3558c2f13df405977111bcb7... f0569cbe09a59336aa2804587cd999ab
20251106_claudesonnet4.5.txt	7026b581d4a77dec04fec3b1bd675544... 9c0a6e4fa1dbed1182aceda31959c5d0
20251106_claudesonnet4.5_translated.txt	4460343edc5b35fcc4131d395e52031e... 7f79e47ed48e80cb672340c048f00314
20251106_gemini2.5pro.txt	bff5f1e5203cb19ca14fe4c42a5b9f24f... 72413c9012ac6dec9097c55aaedfde80
20251106_gemini2.5pro_translated.txt	92bb62b6a8c272c75e0f3202a86d22a8... 5363304a53c4df9b928bbcb4f039f5da
20251108_chatgpt5pro-2.txt	2744aa66e6ea64e365476ad57f5fbba9... 33715f01012487f9166278390437b941
20251108_chatgpt5pro-2_translated.txt	c8ff9752c474eed53163fc39b7108b72... b6800e201fdb2a1df1f7ae7e6fafecf7
20251112_felo.txt	7f693611e23605aa8bbd86d329770a0a... 11cea10ca7ee0f0d640be66992c679cc9
20251112_felo_translated.txt	b618e80af1620e9a1f8fa0f45446a413... 3fbef8352ff0e69310a1d04d20d6fffb0
20251112_genspark.txt	49e3b0f4ec27042ee59fe24567ff6e7b... 96aa95a23bf2cb482bda0d4c36944701
20251112_genspark_translated.txt	19a3427f3205501018e5afddd1aa00bb... 6329811541d2deea915a1d2f62ce349f
20251112_perplexity.txt	6b726e165938ff0364d23395a9c5d098... f33fd3b75943038c7403630d0254b3c
20251112_perplexity_translated.txt	16d740e89777e756faab1c6deaa7ef9... 4419e6252412276fa6e7d007270b2ef7
20251114_gemini2.5proDeepResearch.txt	17a7103c3ae06fe0337eca0b23662270... 4483b1d874cc59a04e5fe3f84d71d7e1
20251114_gemini2.5proDeepResearch_translated.txt	d6b6f84fdb0f11111ea938ac4be850d8... 76ffbc868335b21c086eaf07fd824957
20251115_grok_expert.txt	a5bb7a782cf61c442ed9c82448ac3bcf... fed3e01c3cb2254163e2cd07f45b6a8a
20251115_grok_expert_translated.txt	7ca4e6f932ae4ffb23f68b52d772f553... 601da0851d8d1c25ef4cf4753a173b31
20251119_gemini3Pro_translated.txt	a1c5639450dc8cc0315b3e48d675046b... 129c2253ce6e2814d7d00f06d30818ba

表 1: 実験ログファイルのSHA-256ハッシュ値