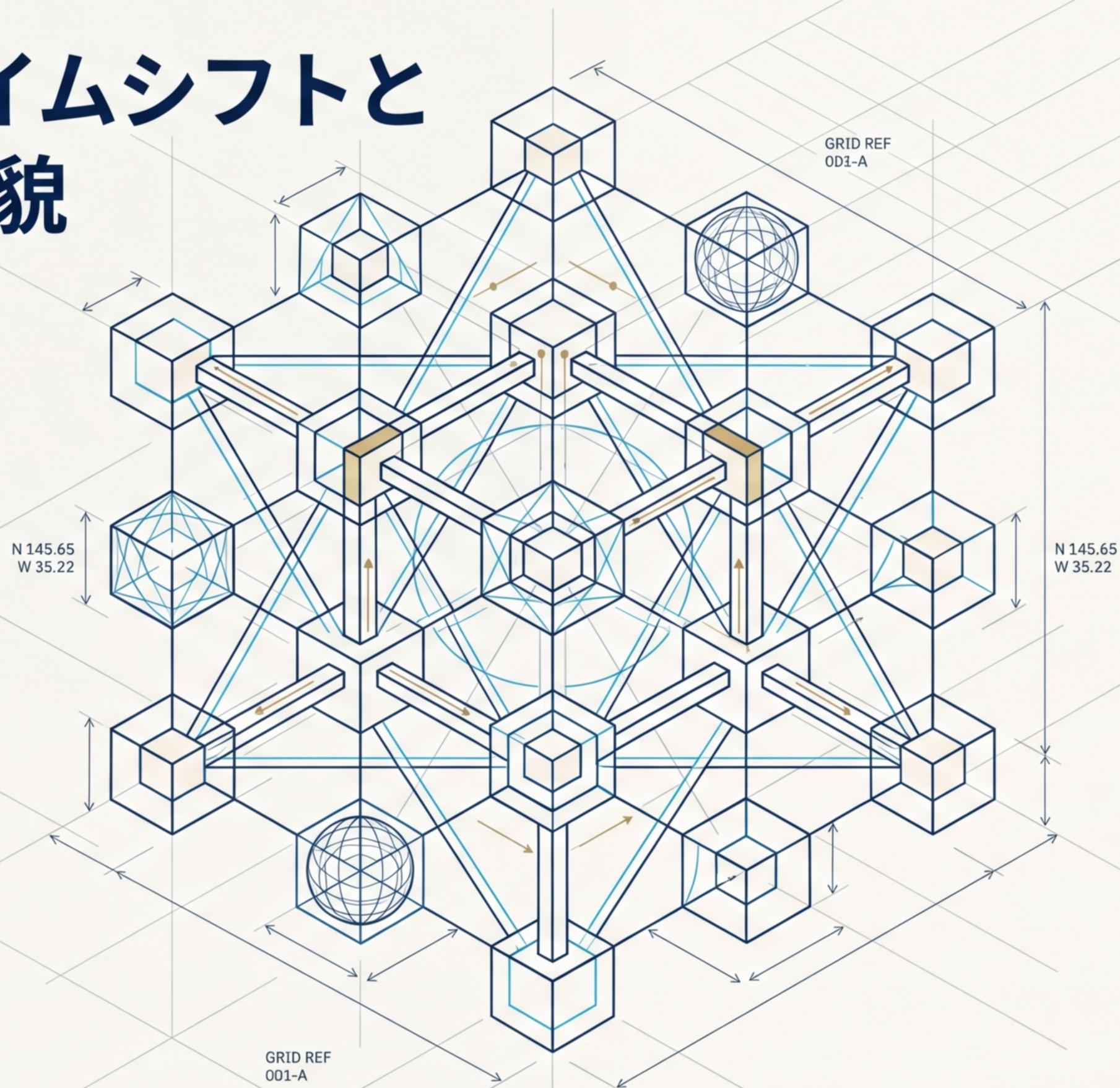


国産生成AIのパラダイムシフトと 「思考エンジン」の全貌

PLaMo 3.0 Prime
アーキテクチャ解析と
グローバル競合比較

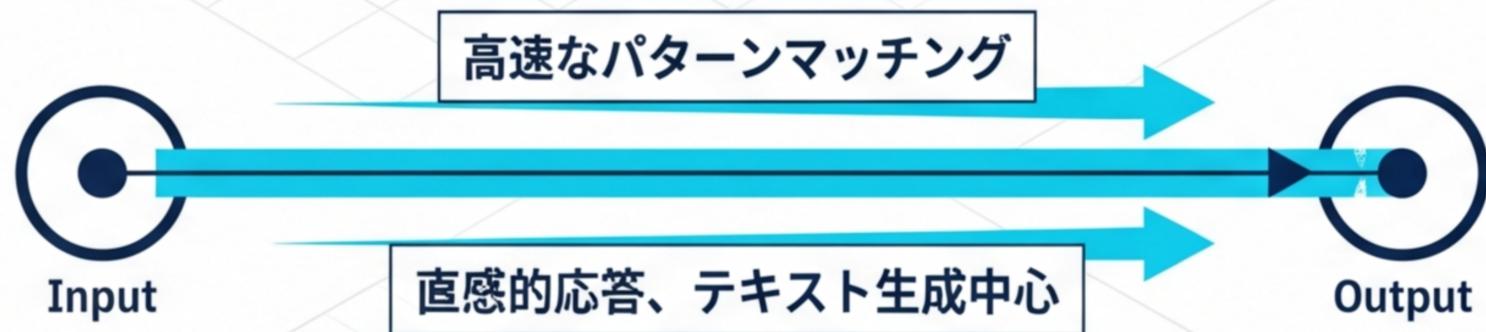


生成AIの劇的進化 — 「System 1」の直感から「System 2」の深層推論へ

従来型LLM (System 1)

COORD: X.024, Y.512

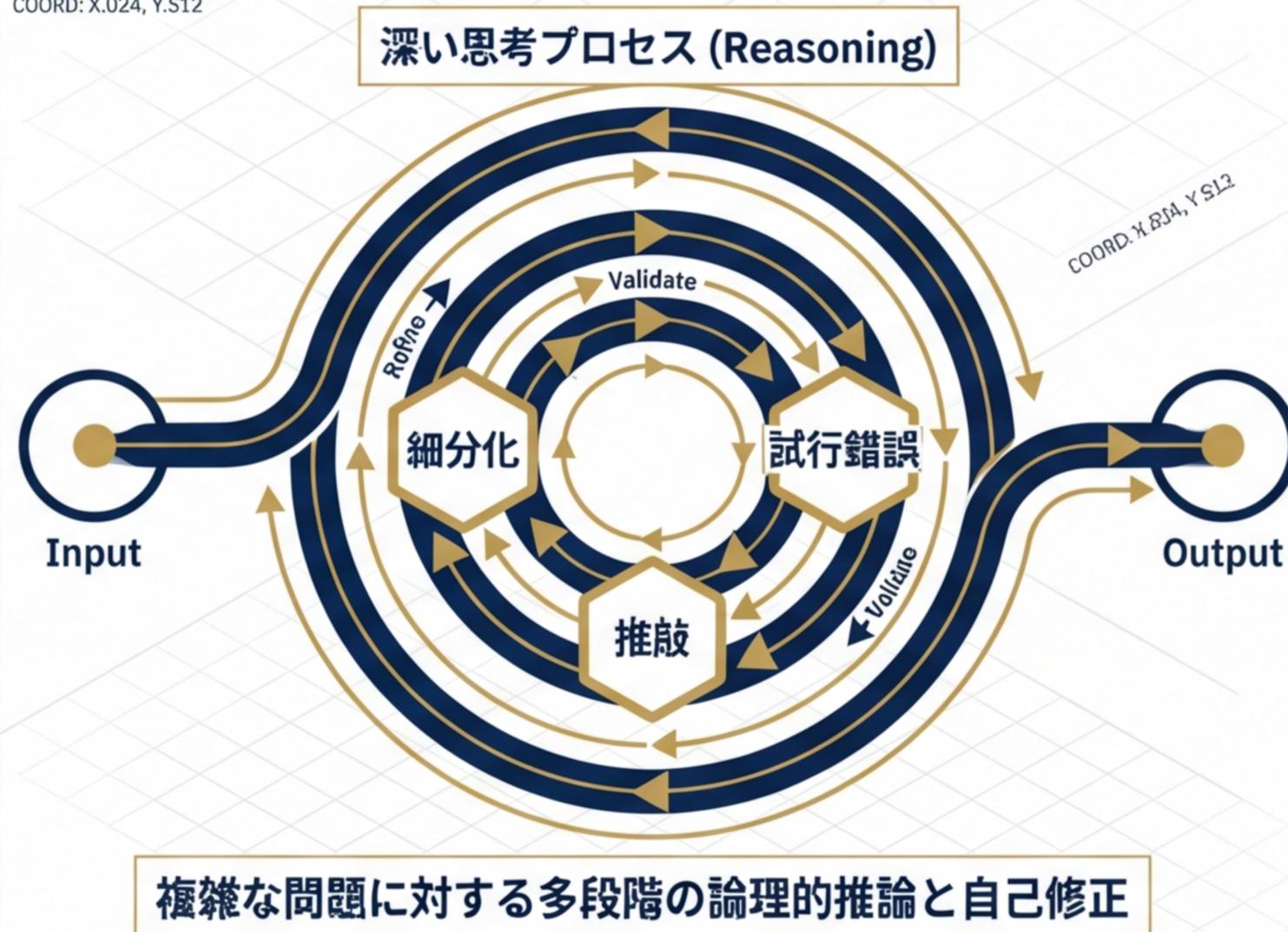
COORD: X.014, Y.S11



推論モデル (System 2)

COORD: X.024, Y.S12

COORD: Y.B14, Y.S13



単なる「テキストの生成」から、複雑な制約条件を解き明かす「論理的で信頼できる思考エンジン」へのパラダイムシフトが完了した。

PLaMo 3.0 Primeの戦略的ポジショニング — フルスクラッチ開発の理由

COORD: X.024, Y.512

Card 1

COORD: X.025, Y.515

ネイティブなReasoning能力の獲得
事前学習段階からのゼロベース構築による、
本質的な推論能力の獲得

COORD: X.026,

COORD: X.024, Y.511

COORD: Y.024, Y.513

Card 2

COORD: X.025, Y.515

日本コンテキストの深層定着
オープンソースの微調整では不可能な、
日本特有の言語的・文化的文脈の
不可逆的な組み込み

Card 3

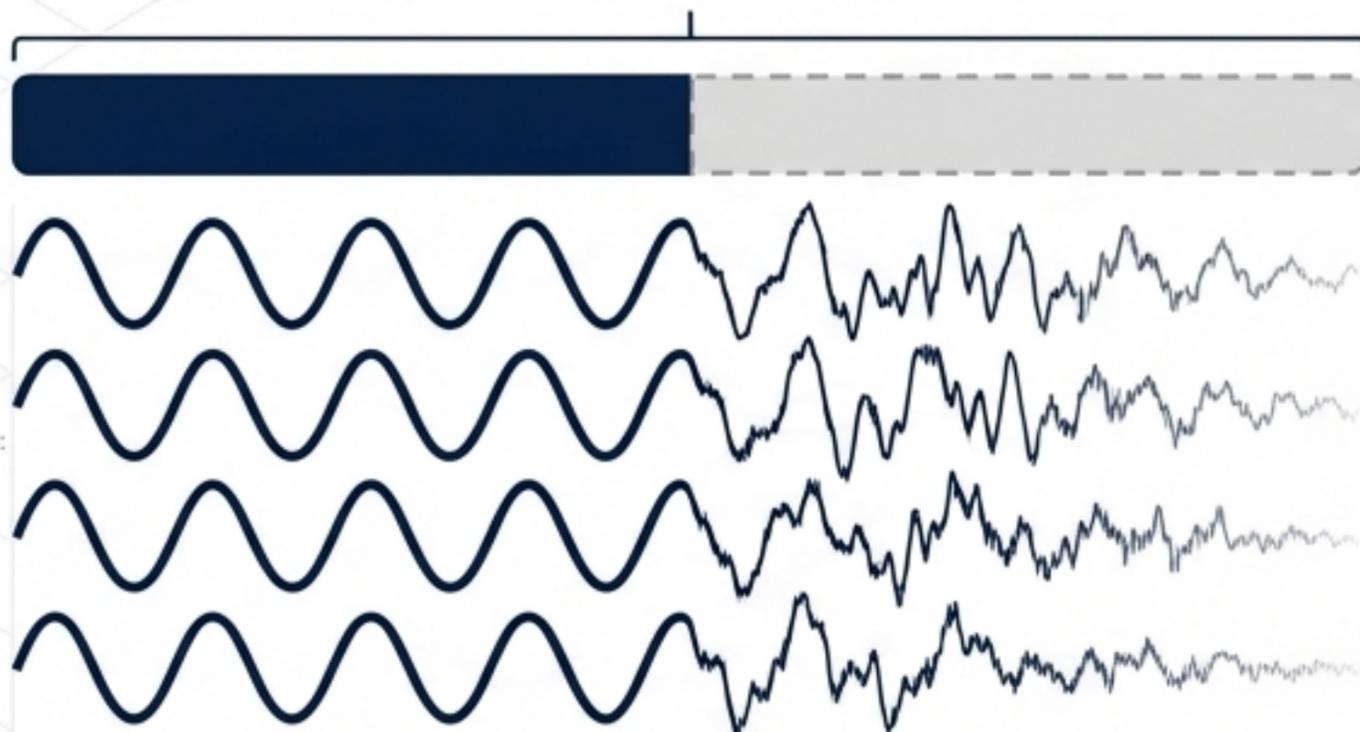
COORD: X.025, Y.513

戦略的意思決定の支援
単なる業務効率化を超え、
企業のコア・コンピタンスに直結する
思考エンジンの提供

既存モデルの「キャッチアップ」を脱し、
独自の思考プロセスを構築するフェーズへの完全移行。

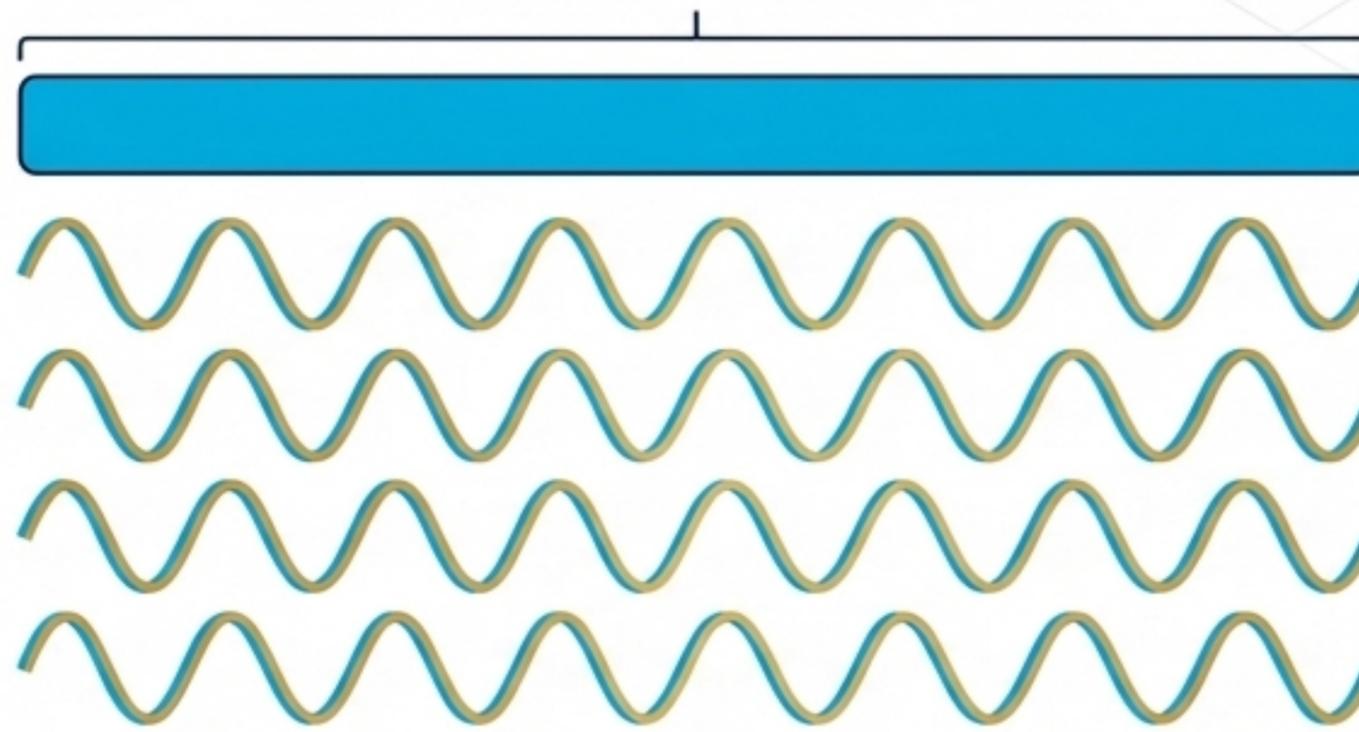
長文脈処理の革新 —— 64Kトークンを支える「YaRN」技術

従来のRoPE

COORD:
X.018

➡ 相対的な位置情報の崩壊と性能劣化

YaRN (コンテキスト拡張)

COORD:
Y.505

➡ 知識を忘却せずアテンション解像度を維持

企業法務における長大な契約書の精査や、巨大なコードベースの解析において、致命的な「文脈の脱落」を完全に防止（最大入力65,536 / 出力20,000トークン）。

「思考エンジン」を構築する3段階の事後学習

事後学習 (Post-Training) パイプライン

COORD: X.025, Y.515

RL (強化学習) 自律的探索と報酬最適化

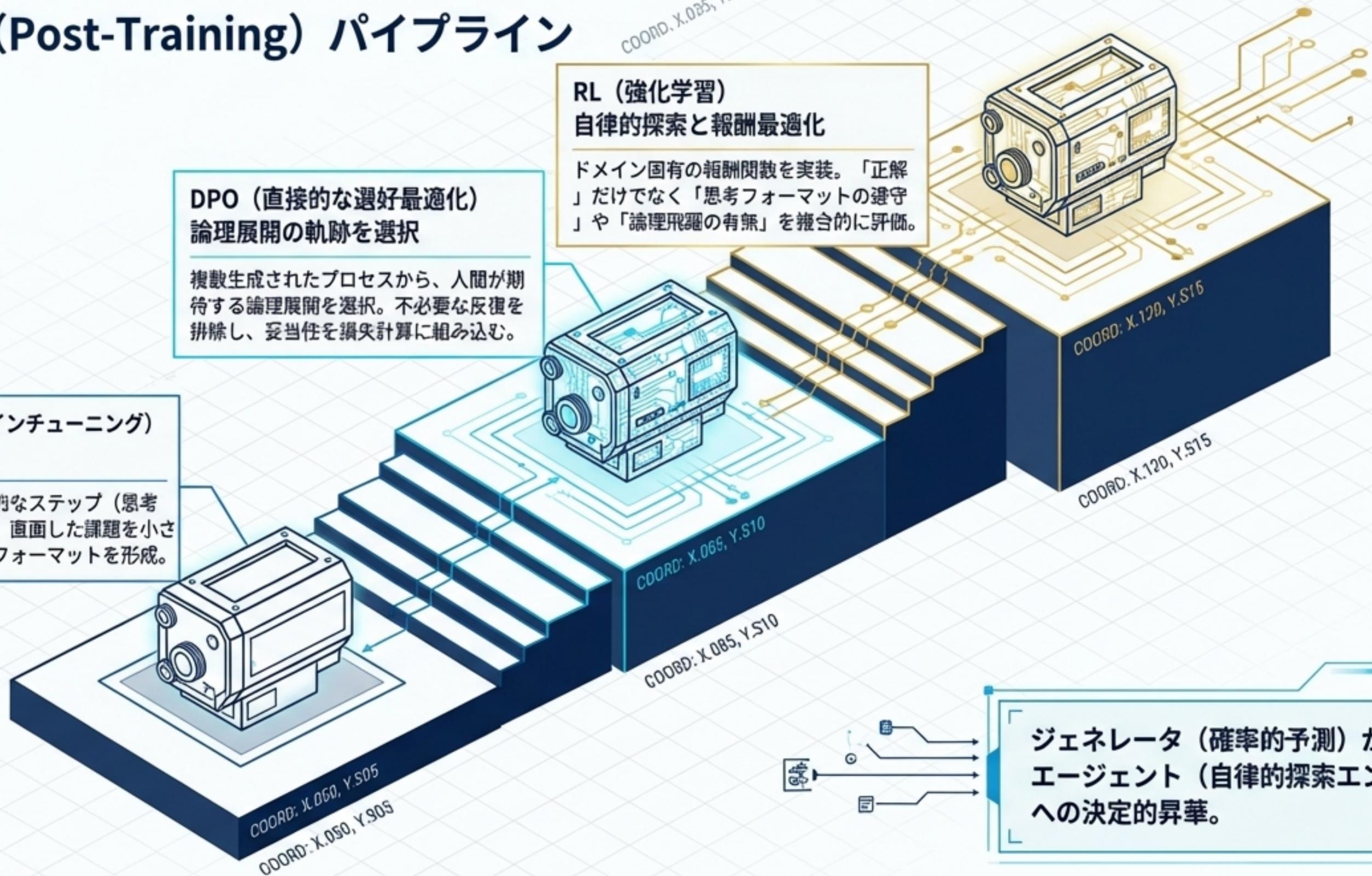
ドメイン固有の報酬関数を実装。「正解」だけでなく「思考フォーマットの遵守」や「論理飛躍の有無」を総合的に評価。

DPO (直接的な選好最適化) 論理展開の軌跡を選択

複数生成されたプロセスから、人間が期待する論理展開を選択。不必要な反復を排除し、妥当性を損失計算に組み込む。

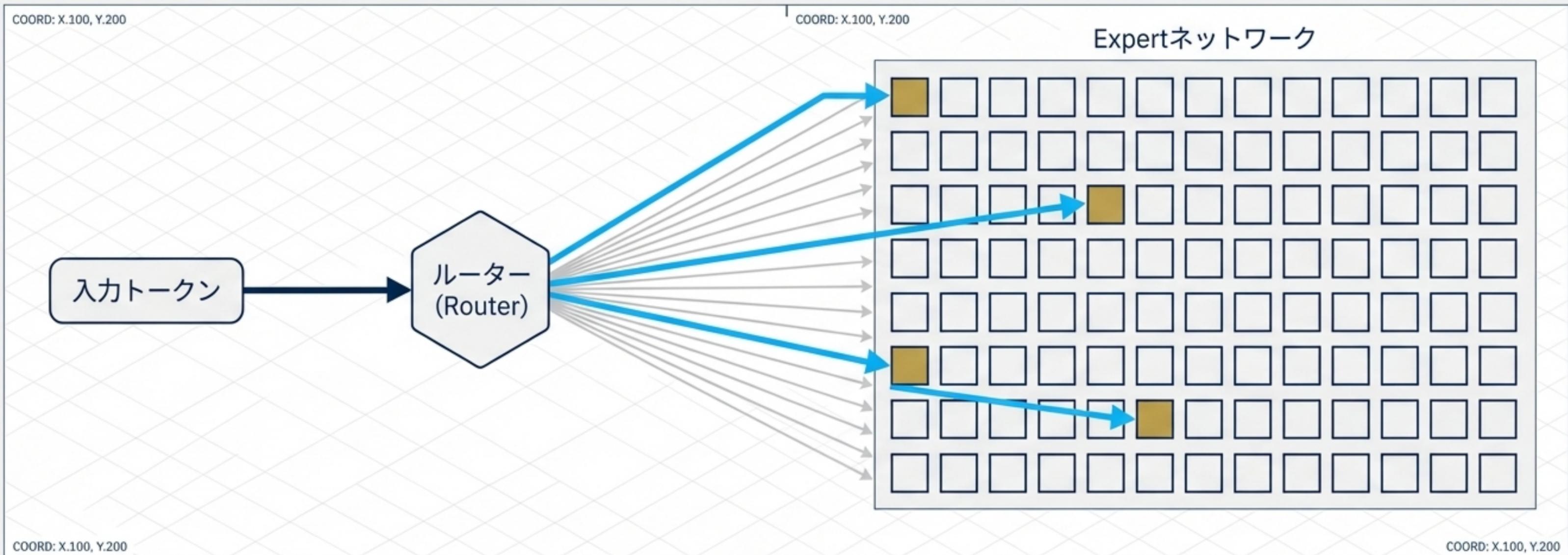
SFT (教師ありファインチューニング) 推論の徹底模倣

結論だけでなく、論理的なステップ (思考過程) を含めて学習。直面した課題を小さな構成要素に分解するフォーマットを形成。



ジェネレータ (確率的予測) から、
エージェント (自律的探索エンジン)
への決定的昇華。

グローバル・トレンド — 熾烈を極める「MoE」アーキテクチャの台頭



MoE (Mixture-of-Experts)

入力トークンごとに少数の「Expert (専門家)」サブネットワークのみを動的に稼働 (ルーティング) させるスパース (疎) な設計。

The Dilemma Solved

「巨大な知識量 (総パラメータ数)」を維持しながら、「推論コスト (アクティブパラメータ数)」を劇的に抑制するためのグローバル標準アプローチ。

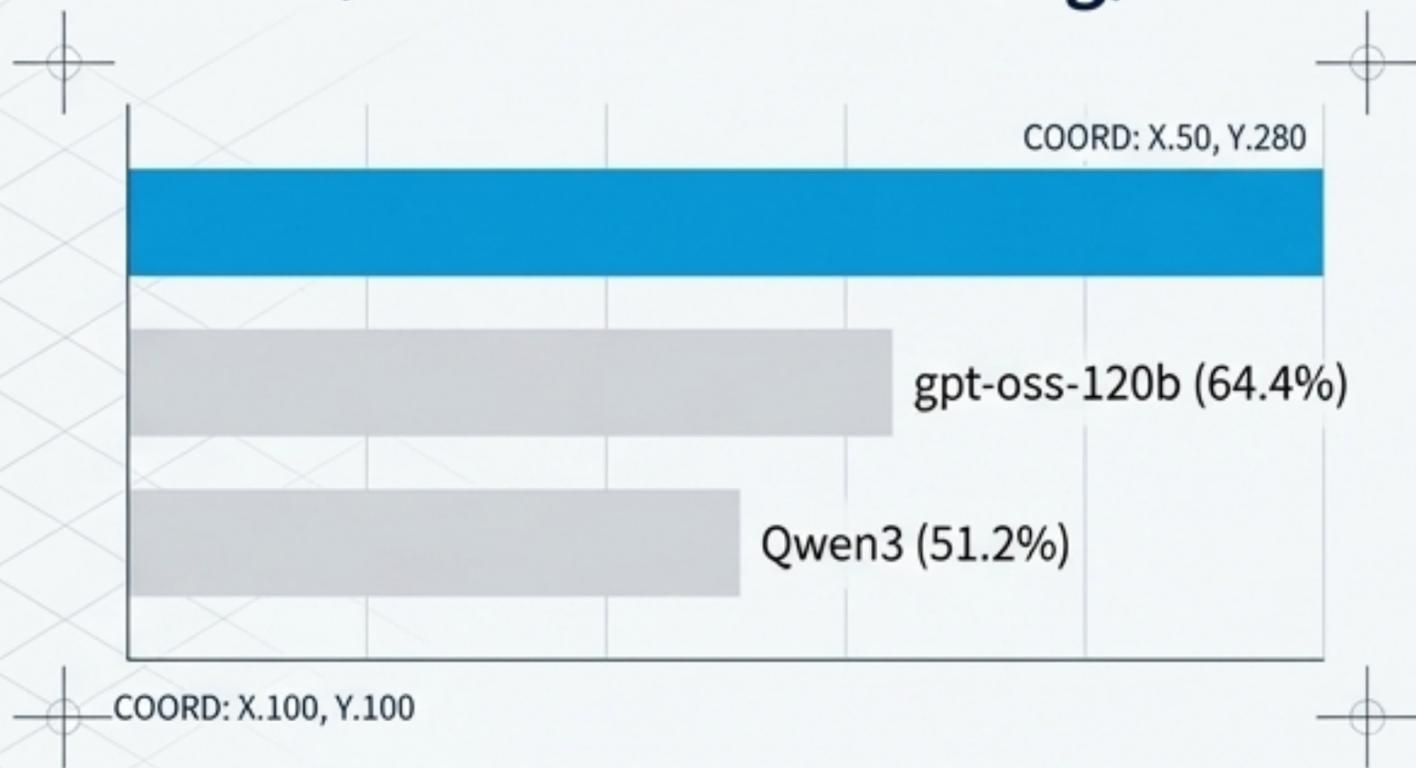
【比較マトリクス】グローバル推論モデルの「効率vsスケール」競争

指標	gpt-oss-120b	Qwen3-235B-A22B-Thinking-2507
総パラメータ	117B	235B
アクティブパラメータ	5.1B (極端なスパース設計)	22B (高リソース要求)
スパース比率	約 4.3%	約 9.4%
API 出力コスト (1M token)	\$0.19 (圧倒的低コスト)	\$1.50 (gpt-ossの約8倍)
ハードウェア要件	単一80GB GPU稼働可能	複数GPU必須

gpt-ossは「効率の極致」、Qwen3は「スケールの極致」を追求。この両極端なグローバル市場において、PLaMo 3.0 Primeはいかに独自の価値を証明するのか？

【ベンチマーク評価①】 実務応用におけるPLaMo 3.0 Primeの圧倒的優位性

指示追従の精緻さ (Instruction Following)



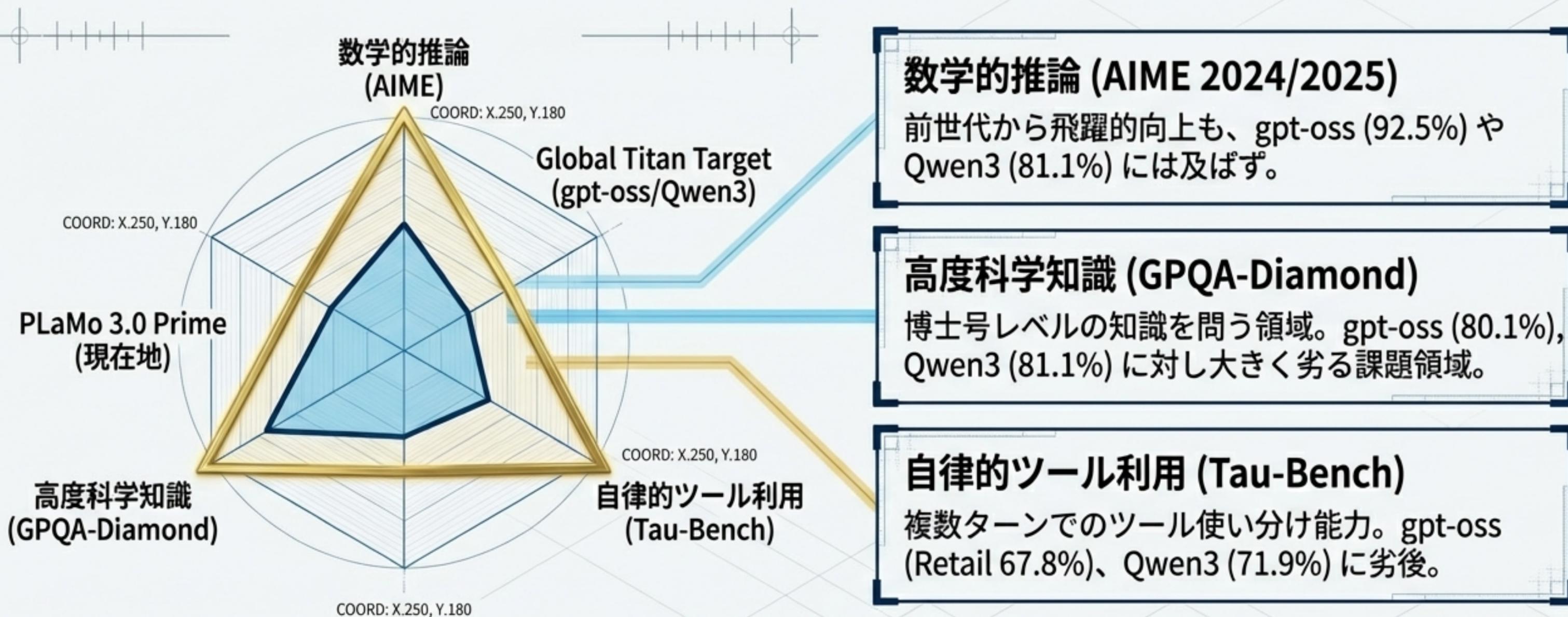
IFBench (英語) / JFBench (日本語) 共に PLaMoが最上位を獲得。複雑な制約事項を遵守する能力において、プロンプトエンジニアリングの要求に最も忠実に応える。

日本コンテキストの完全把握



Japanese MT-Benchにて「性能の天井」に到達。複雑な敬語表現や、文脈に依存した曖昧なニュアンスの処理において、グローバルモデルと同等以上の高度な出力を実現。

【ベンチマーク評価②】 深層推論とエージェント機能における「現在地」



数学的推論 (AIME 2024/2025)

前世代から飛躍的向上も、gpt-oss (92.5%) や Qwen3 (81.1%) には及ばず。

高度科学知識 (GPQA-Diamond)

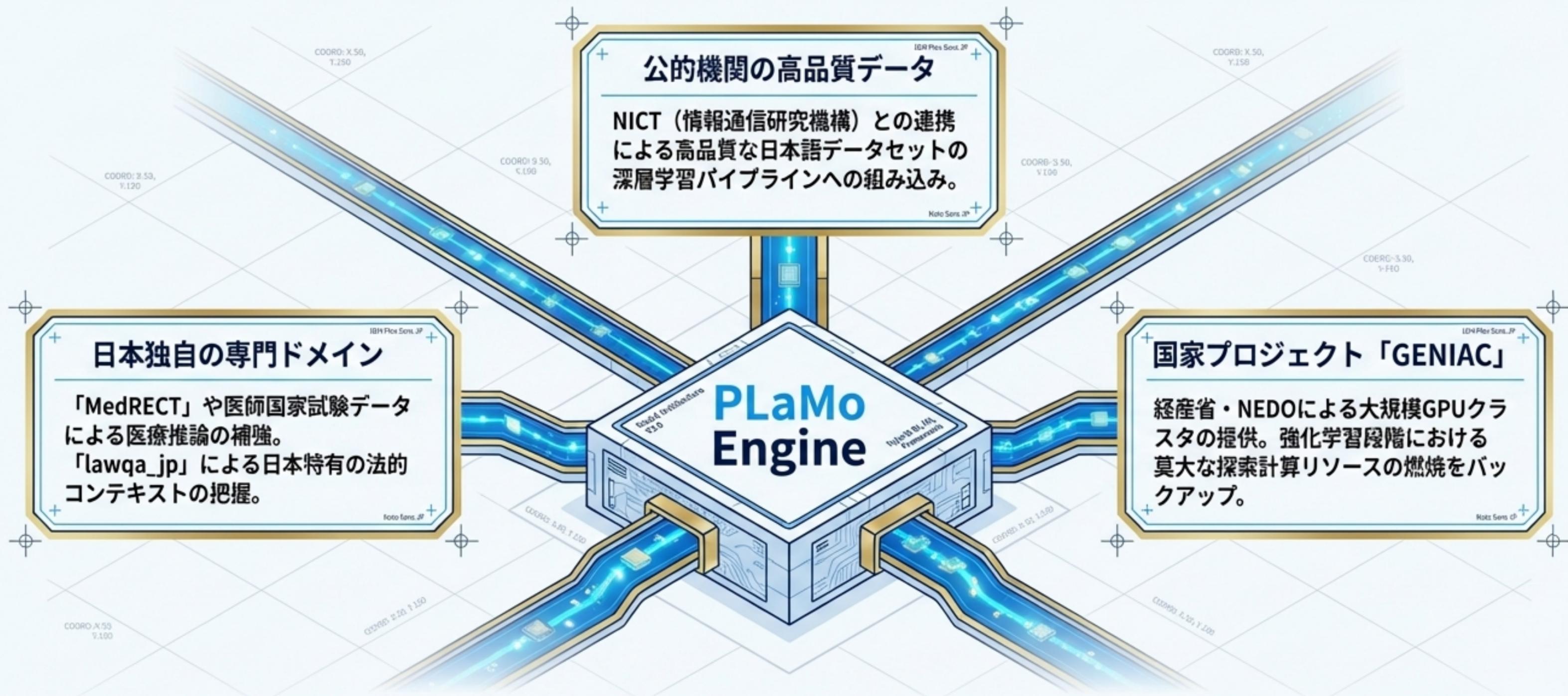
博士号レベルの知識を問う領域。gpt-oss (80.1%), Qwen3 (81.1%) に対し大きく劣る課題領域。

自律的ツール利用 (Tau-Bench)

複数ターンでのツール使い分け能力。gpt-oss (Retail 67.8%)、Qwen3 (71.9%) に劣後。

学習データボリュームの絶対差と、RLスケールの差が要因。PFNはこれらを隠さず、β版を通じた最優先の改善項目として透明性高く提示している。

弱点を凌駕する学習データ戦略と「GENIAC」による国家規模の支援



外資系モデルがカバーしきれない「日本の機微」を理解する推論能力の獲得。



COORD: X.300,
Y.150

エンタープライズ導入における究極の価値 「データ主権」の完全確立

セキュリティ/主権の強度

パブリックAPI



COORD: X.300,
Y.150

VPC環境



完全オンプレミス

COORD:
X.300, Y.150
Y.130



- 外資系クラウドへのデータ送信 / 経済安全保障リスク

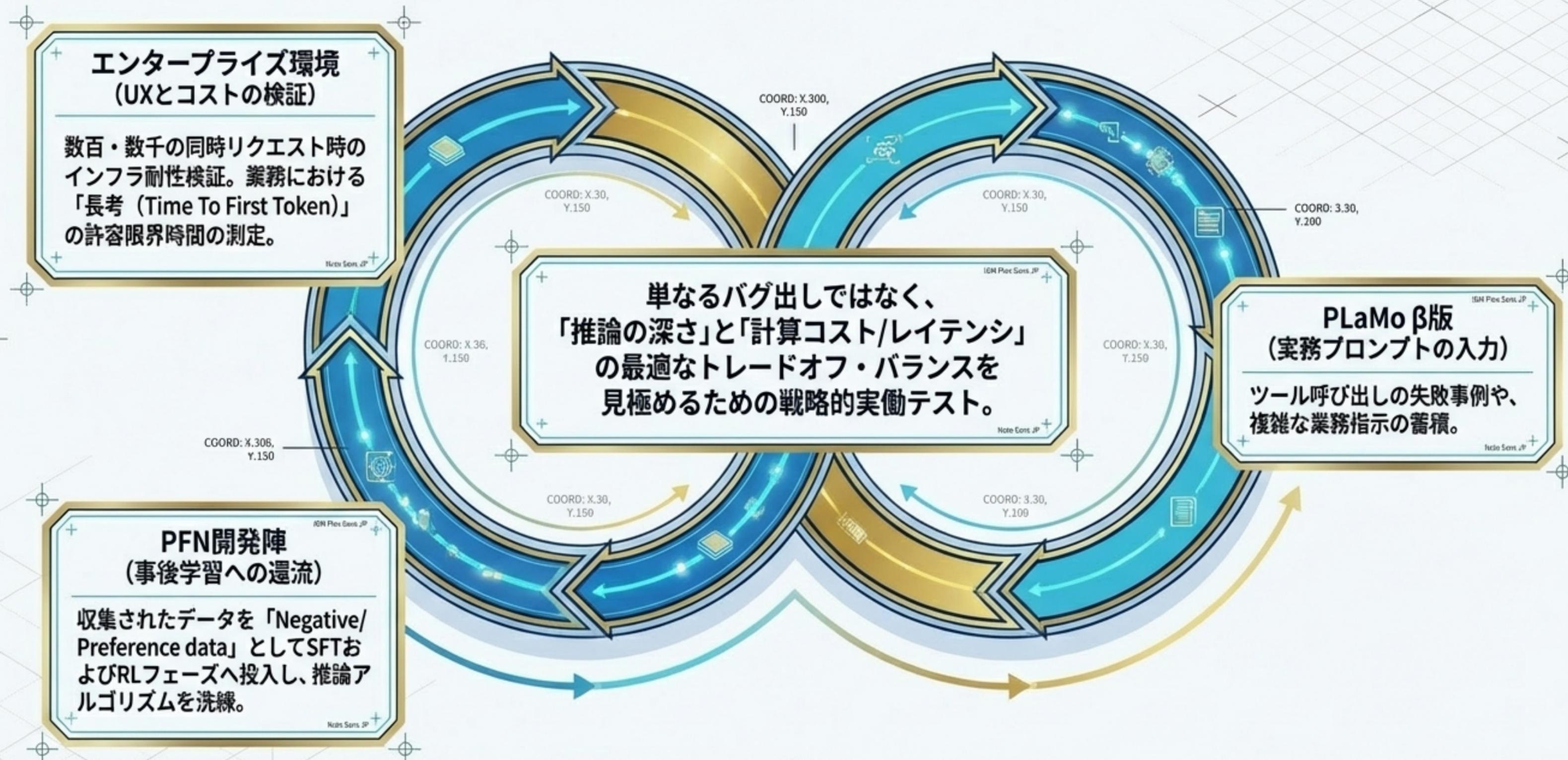
COORD: X.306,
X.150, Y.150

- Amazon Bedrock / Snowflakeへのセキュアな統合

- 自社データセンター内の物理サーバー稼働

M&Aの財務分析、未公開特許の探索、新薬開発など「企業のコア・コンピタンス」に関わる機密データを、外部通信を完全に遮断した状態で、日本特有の商習慣に精通した思考エンジンに処理させる圧倒的優位性。

β版モニタープログラムの真の目的 — 実運用データによるエコサイクルの始動



結論 (Executive Summary) — デジタル主権を築く「思考エンジン」の未来

1

推論の完全国産化

オープンモデルの「キャッチアップ」時代は終焉。自国の文化・ビジネスロジックに根ざした独自の思考プロセスをフルスクラッチで構築。

2

実用と主権の両立

特定ドメインにおける圧倒的な指示追従性と、完全オンプレミス展開によるデータ主権（セキュリティ）の確立。

3

共に育てるAIエコシステム

2026年6月の商用版リリースに向け、日本のエンタープライズ企業との共創を通じて推論の精度と運用効率を極限まで高める。

PLaMo 3.0 Primeは、日本企業に対してAI戦略の新たな、そして不可欠な選択肢をもたらす。