

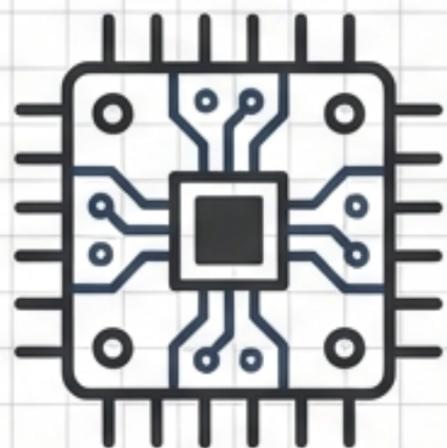
Executive Briefing

# 初の“長考”できる国産 フルスクラッチLLM 「PLaMo 3.0 Prime」 深掘り分析

Reasoning能力の真価と、企業導入  
に向けた技術・運用ロードマップ

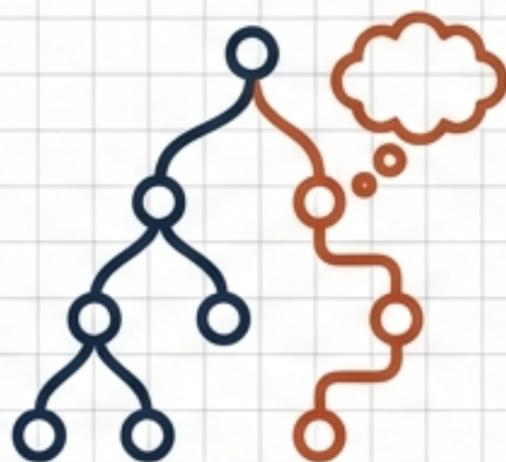
# Executive Summary

## 1. 開発の現在地



- PFNがゼロベースで再構築したフラッグシップ。
- 最大64Kトークンの長文コンテキストを処理可能な国産フルスクラッチLLM。

## 2. 核心的イノベーション



- 回答だけでなく「思考過程」を出力する“**Long-form Reasoning**”（長考）機能を実装。
- 事後学習（SFT/DPO/RL）において、**思考プロセス自体**に損失計算・報酬関数を適用。

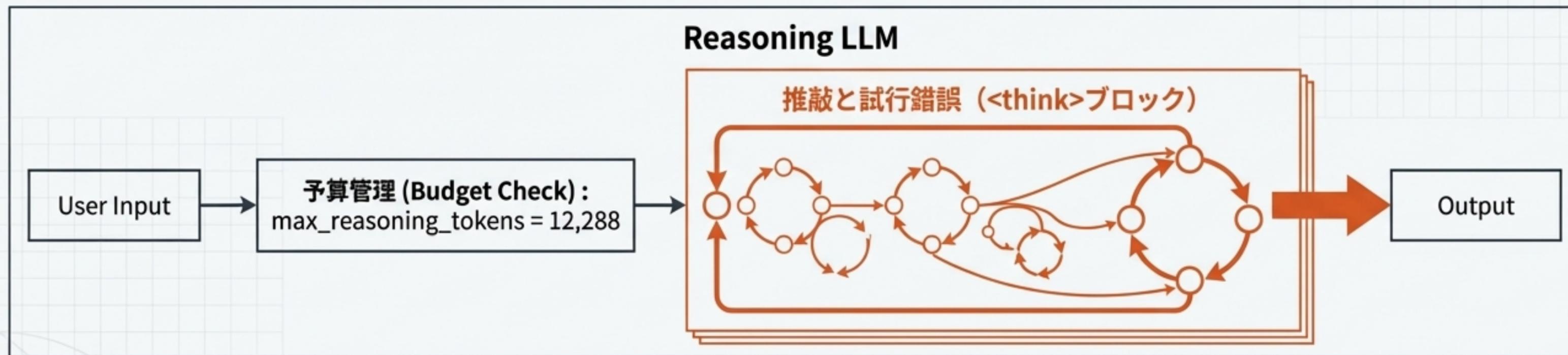
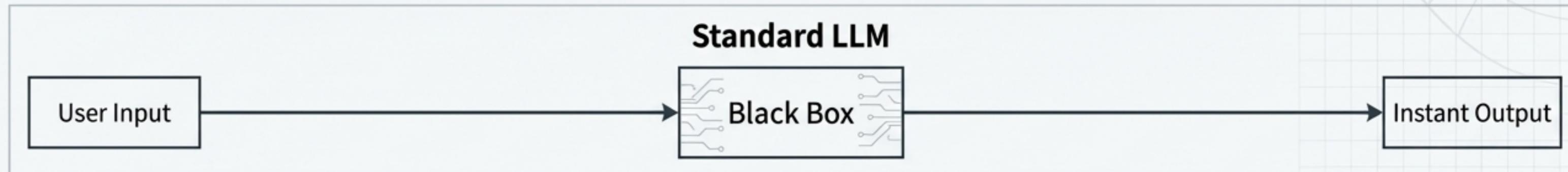
## 3. 企業導入のリアリティ



- 日本語・英語の対話や法務・医療タスクでグローバル最前線に匹敵。
- 推論時の計算コストやエージェント機能の成熟度に課題。実運用環境での $\beta$  モニター検証が必須。

# 「瞬発的な回答」から「予算化された思考プロセス」への転換

長考 (Long-form reasoning) とは何か？



## 探索的推敲

小さなステップに分割し、**推敲**を繰り返す過程を出力。

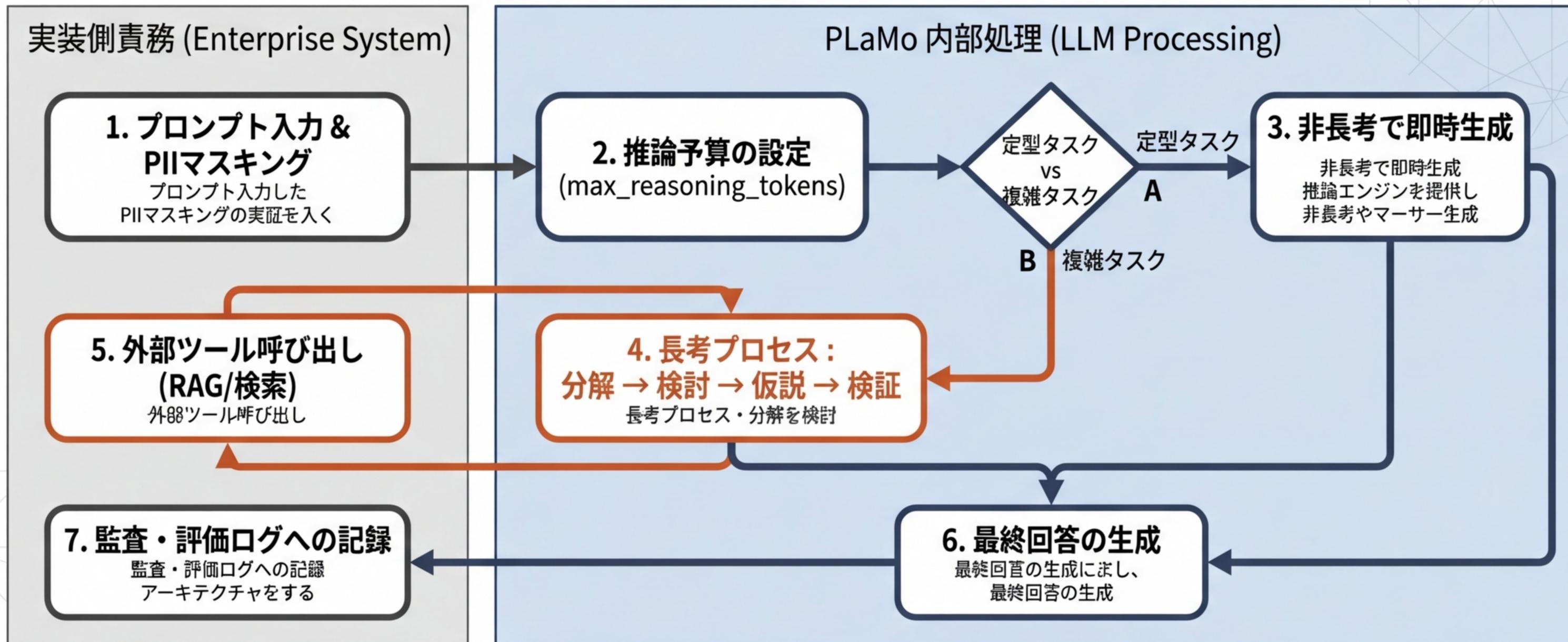
## 思考の最適化

SFT/DPOで「最終回答」だけでなく「**思考過程**」自体を評価。

## 推論の予算化

計算リソースを「**思考用トークンの予算**」として明示的に制御。

# 企業システムにおける“長考”の処理フローと責任分界点



PLaMoは高度な推論エンジンを提供するが、ツール実行やデータマスキングは企業側のアーキテクチャに依存する。

# PLaMo 3.0 Prime β版の技術プロフィール（解剖図）

## 公表済み仕様（Confirmed）

### コンテキスト長

入力 64K / 出力 最大 20K  
(YaRN採用)

### 学習パイプライン

事前学習 → 継続事前学習  
→ SFT → DPO → RL

### 学習データ

独自データ + 医療特化  
+ NICT日本語データ

## アーキテクチャの推定（Inferred）

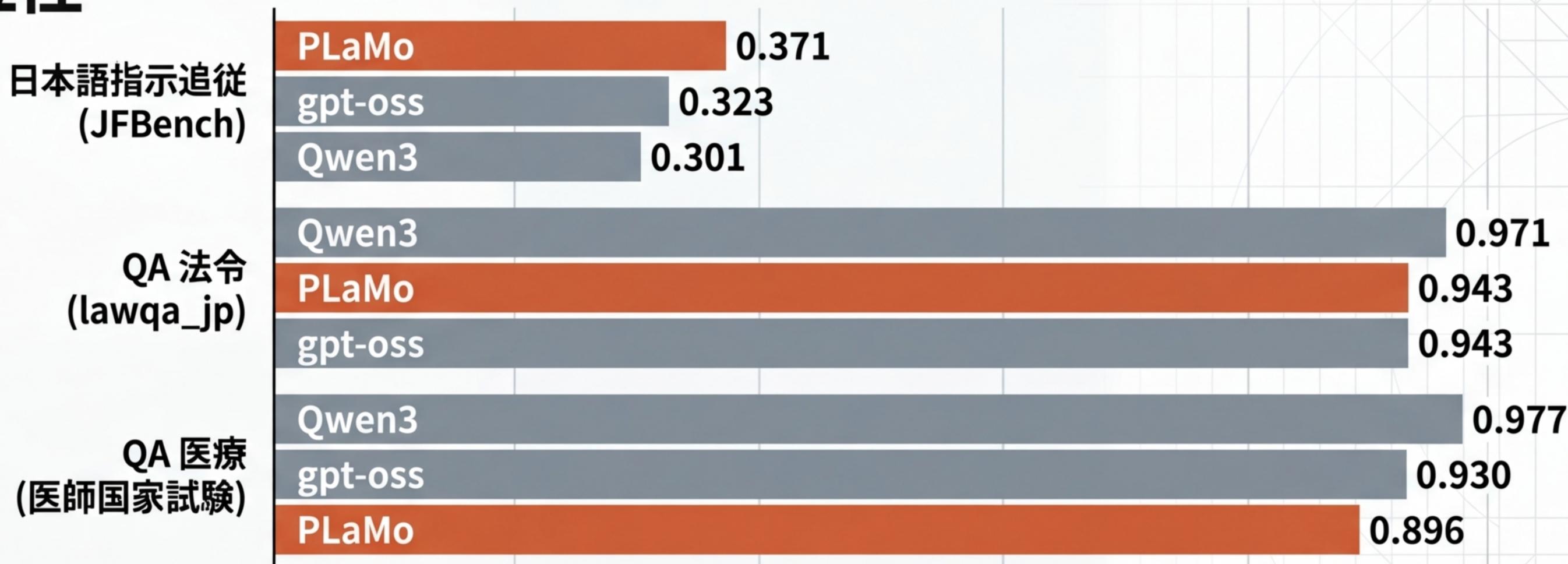
### 高効率設計

PLaMo 3系の傾向から、Sliding Window Attention等を採用しKVキャッシュ消費を抑制している可能性。

### モデル構造

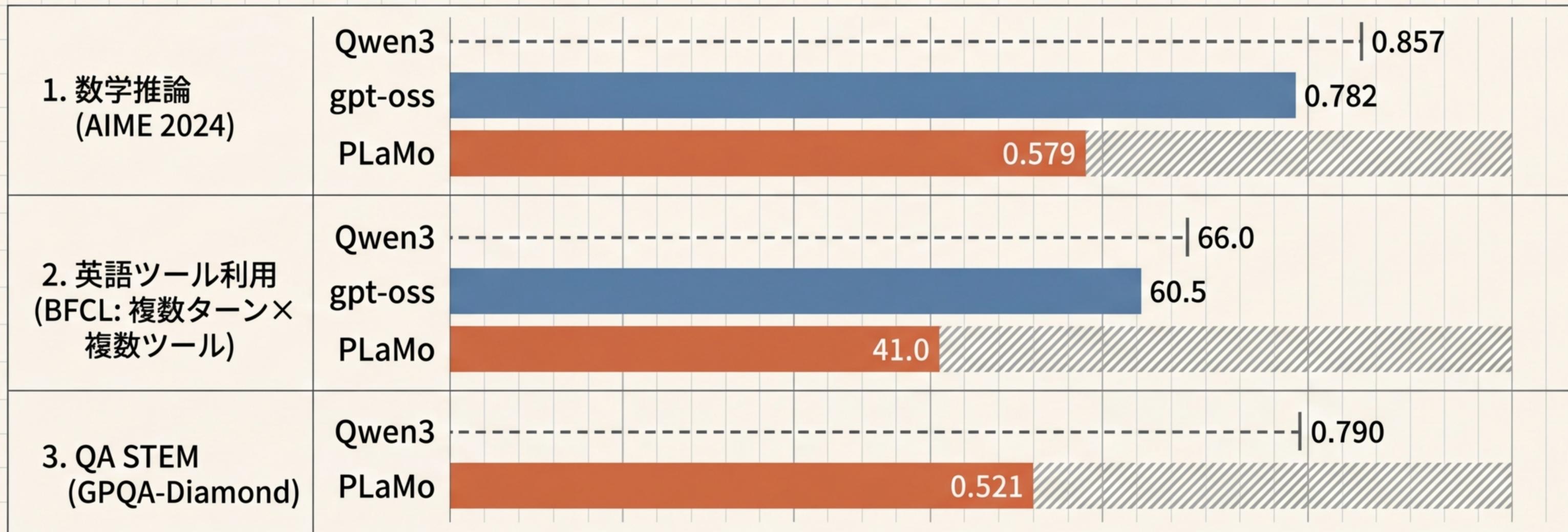
DenseまたはMoEかは未公開。競合トレンドから大規模かつ高効率な活性パラメータ設計と予想。

# ベンチマーク分析 (1/2) : 国内実務・専門領域における優位性



ドキュメントヘビーな国内実務（規程チェック、要件定義、医療・法務QA）において、グローバルなオープンウェイト巨人に肉薄する性能を実証。日本語のニュアンスとコンプライアンス要件が絡む業務に最適化されている。

# ベンチマーク分析 (2/2) : エージェント化と数理推論における現在地



外部環境からのフィードバックを伴う「動的ツール選択」や「高度な数理ロジック」においては、グローバル水準との差が存在する。“ツール連携で長考を拡張する”自律型エージェント領域は、今後のアップデートの焦点となる。

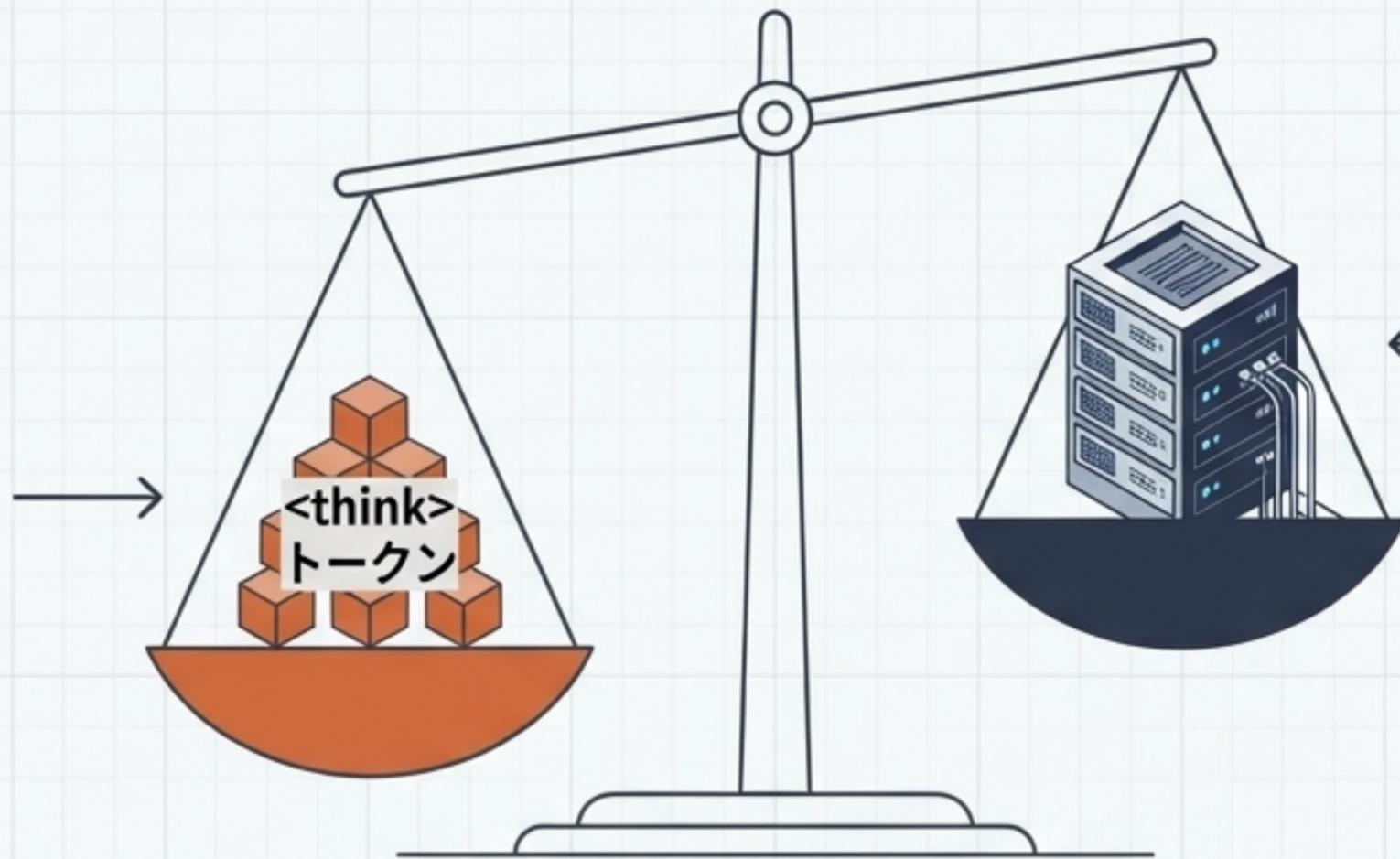
# 競合マトリクス：グローバル・オープンモデルとの構造比較

	PLaMo 3.0 Prime	Qwen3-235B	gpt-oss-120b
開発形態	フルスクラッチ国産	オープンウェイト	オープンウェイト
アーキテクチャ規模	未公開	MoE 235B/活性22B	MoE 117B/活性5.1B
コンテキスト長（入力）	64K	131K YaRN	128K
出力上限	最大20K	32K	可変
長考の推論制御	max_reasoning_tokens=12,288	thinking mode	reasoning effort (low/med/high)

PLaMo 3.0 Primeは、推論トークン量を明示的に制御可能な「国産フルスクラッチ」として独自のポジションを確立。

# “思考コスト”の可視化：推論計算量とインフラ要件のトレードオフ

**Reasoning Depth:**  
max\_reasoning\_tokens=12,288 の消費。  
長考するほど出力生成前の内部トークン量が激増する。



**Infrastructure Cost:**  
VRAM推定470GB級の重みメモリと、長系列トークン生成による推論レイテンシの増加。

## The Language Factor (英語思考の課題)

現状、学習データの影響で「思考プロセスが主に英語」で行われている。英語での長考はトークン消費効率を悪化させるため、正式版に向け「思考の日本語化」によるトークン節約を開発中。

長考は無料ではない。精度とレイテンシのバランスを設計する「推論予算 (Compute Budget) の管理」が企業側の必須スキルとなる。

# アーキテクチャに適合するターゲット・ユースケース

## 長文ドキュメントの 統合解析

強み: 64Kコンテキスト ×  
高度な日本語指示追従。

適用例: 大量の社内規程、  
契約書群、要件定義書を  
跨いだ整合性チェックと  
論点抽出。

## 高信頼性が要求 される専門的QA

強み: 医療・法令ベンチマー  
クにおける圧倒的スコア。

適用例: RAGと連携した法  
務部門の一次ドラフト作  
成、医療情報の専門的検  
索補助。

## 複数条件が絡む 意思決定支援

強み: 段階的推論による  
思考プロセスの可視化。

適用例: 複雑な申請審査、  
社内コンプライアンス適  
合チェック、根拠提示が  
必要なサポート対応。

# デプロイメント・タイムラインとβモニターの真の目的

2025年11月

PLaMo 3  
小規模事前学  
習検証。高効率  
アーキテクチャ  
の基礎検証。

2026年3月  
[現在地]

PLaMo 3.0  
Prime β版  
モニター  
募集開始。

## なぜβテストが必要か？

モニターの最大の狙いは「精度の検証」だけではない。長考モデル特有の「トラフィック量」「実運用時のレイテンシ」「並列処理時の品質ブレ」を現場の負荷環境でストレステストし、商用インフラの安定性を担保することにある。

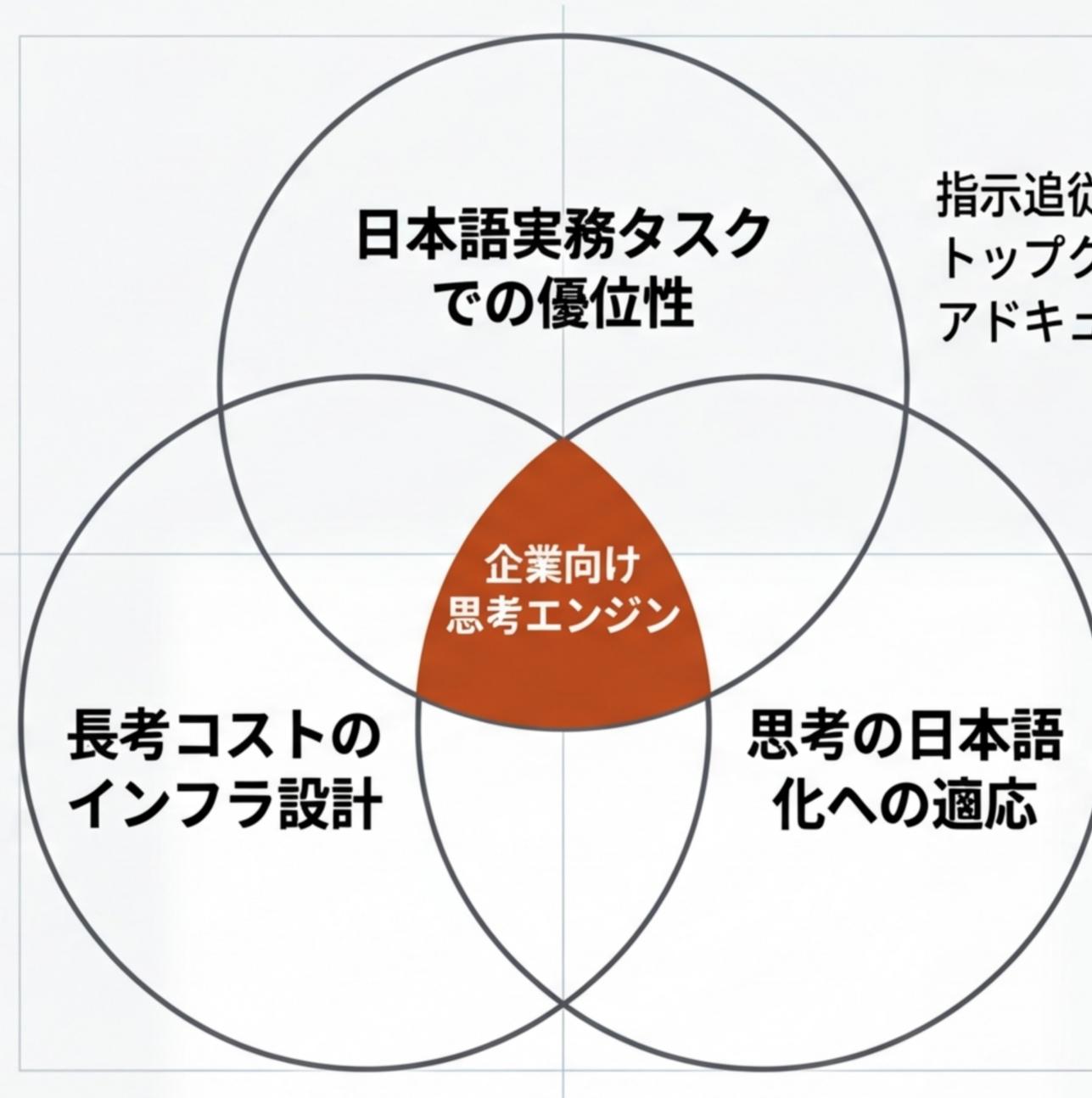
2026年6  
月中旬

商用版  
正式提供  
開始予定。

# 推論モデル特有のリスクとAIガバナンス実装要件

<h2>1. 安全性（幻覚・過信リスク）</h2> <p>Risk: もっともらしい「思考プロセス」が誤った結論を正当化する。</p> <p>↳ Mitigation: 根拠資料（RAG）の必須化、<b>人手レビュー</b>（Human-in-the-loop）による最終判断。</p>	<h2>2. プライバシー（データ混入）</h2> <p>Risk: 長文コンテキスト入力による個人情報の意図せぬ混入。</p> <p>↳ Mitigation: プロンプト入力段階での<b>PIIマスキング</b>と、<b>監査可能なデータフロー</b>の確保。</p>
<h2>3. 著作権・知財</h2> <p>Risk: 文化庁ガイドライン等に抵触する類似出力のリスク。</p> <p>↳ Mitigation: 生成物の類似性を下げる<b>出力フィルタリング</b>と、社内法務ワークフローの統合。</p>	<h2>4. バイアスと表現癖</h2> <p>Risk: 「英語中心の思考」や翻訳データ由来のニュアンスの偏り。</p> <p>↳ Mitigation: 業務KPIに直結する<b>ドメイン特化のレッドチーミング</b>とバイアス評価の実施。</p>

# Strategic Synthesis : PLaMo 3.0 Primeを事業価値に転化する3つの条件



指示追従や長文処理における国内トップクラスの精度を、自社のコアドキュメントで活かし切る。

高騰しがちな推論コストに対し、「推論予算の制御」と「必要時のみ長考を回すルーティング」を実装する。

正式版に向けた日本語推論能力の向上を取り込み、トークン消費効率の劇的な改善を達成する。

PLaMo 3.0 Primeは単なるチャットボットではない。推論予算を投資して確実なロジックを組み上げる「企業向け思考エンジン」としてのアーキテクチャ再設計が、導入成功の鍵を握る。