

2026/4/24 財務省財務総研

# 生成AIの進化と未来

## - 国産AIとフィジカル AI

Preferred Networks 共同創業者 代表取締役社長  
Matlantis 代表取締役社長  
岡野原 大輔

# 自己紹介：岡野原 大輔

Preferred Networks 共同創業者 代表取締役社長

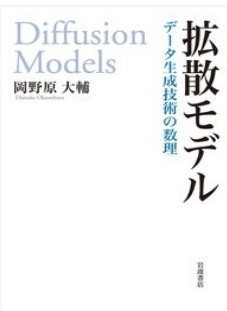
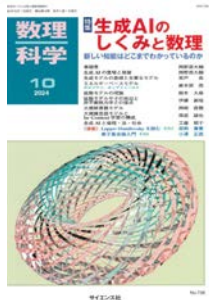
Matlantis 代表取締役社長

X (Twitter): @hillbig

AI事業者ガイドライン 検討会委員

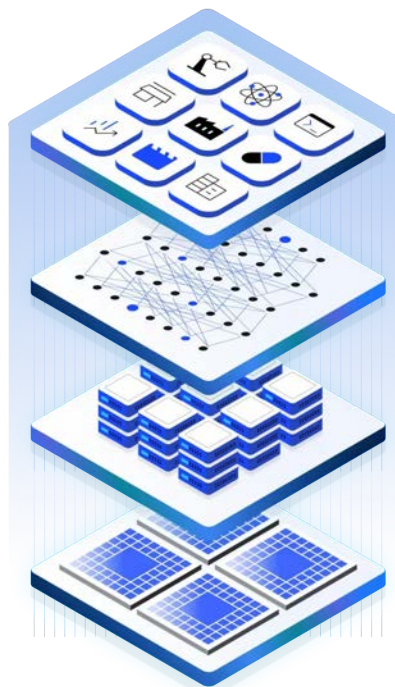


## 著書



# Preferred Networksの事業:

PFNは、チップ、計算基盤、生成AI基盤モデル、ソリューション・製品まで、AI技術のバリューチェーンを垂直統合し、ソフトウェアとハードウェアを高度に融合した事業を進めています。



## AIプロダクト・ソリューション

様々な産業向けのAIソリューション・製品

PreferredAI

MiseMise

MATLANTIS

PFN 3D Scan

Preferred Networks  
Visual Inspection

kachaka

## 生成AI基盤モデル

PLAMO

大規模言語モデル

PFP

物質のエネルギー計算モデル

## 計算基盤



GPUクラスタ



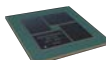
MN-3  
(MN-Coreクラスタ)

PFCP

MN-Core 2を計算資源とした  
クラウドサービス

## AI半導体

MN-Core



MN-Core



MN-Core 2



MN-Core L1000  
(2027年提供予定)



次世代

# AIはどう作られるか

これは何の絵でしょう？



これは何の絵でしょう？ → 馬の絵 (2013年の生成AI)



# これは何の絵でしょう？ 馬(2026年の生成AI)



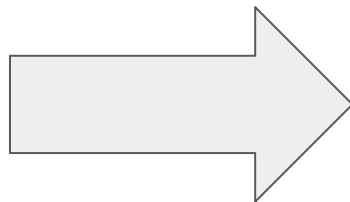
# 生成能力の驚異的な進化速度

生成はわずか10年間で驚異的に進化した。様々な生成（言語、画像、ロボット制御）が統一的に扱われ急激に進化した

(生成データ)(学習データ)



動物の画像生成  
2013年



2026年の画像生成

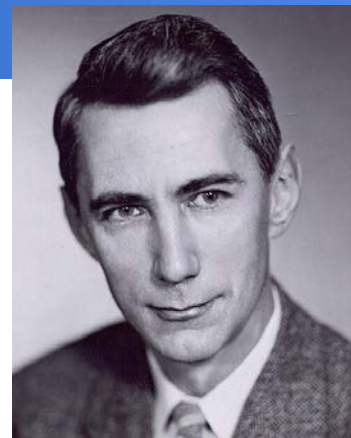
# 言語モデル [Shannon 1948]

これまでの単語に後続する単語を順に予測するモデル

私は 父と 一緒に \_\_\_\_\_

「病院」? 「遊ぶ」? 「過ごした」?

ここにどの単語が出現しそうかを予測する



Claude Shannon

シャノンが1948年に情報理論を確立した最初の論文でも、次の文字を予測するモデルを作ることのできるような文字列が生成されるのかを実験している

「OCRO HLI RGWR MMIEL」 (1文字頻度情報を使った場合)

「IN NO IST LAT WHY CRA」 (3文字頻度情報まで使った場合)

# 言語モデルは単にオウム返しをしているだけか？ 違う！

推理小説を全部読み終わった後に

「そうかわかった！ 犯人は ...」

の次の単語を予測できるなら、そのAIは推理小説を全部理解し、推定までできていることになる。

次の単語予測を「知能を鍛える練習問題」とし、解けるようにする中で、結果として言語やその背後の概念を理解する

**大規模データを使った知能の学習手法が確立**

# なぜAIはこんなに急速に進化したのか？

## 1. 計算機が急速に速くなり続けている

AIを作るのに使われる計算力は毎年4.6倍

## 2. AIを作る(学習)ときに使われるデータの指数的增加

現在は数十兆文字を学習させている

## 3. たくさんの人が研究している(年間数十万報)

毎日数百本の研究成果・SWがリリース

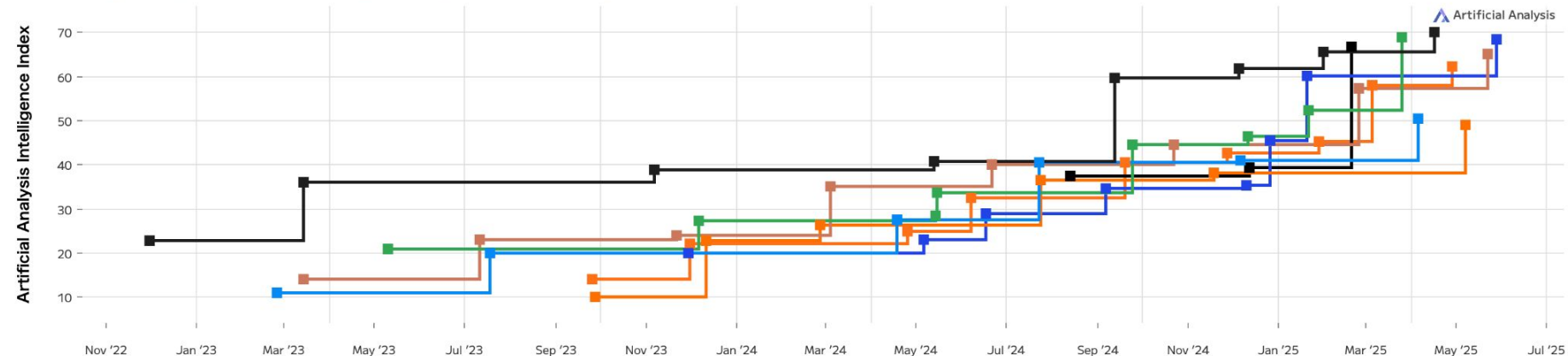
# AIの進化

# AIは淡々と少しずつ賢くなっている

## Frontier Language Model Intelligence, Over Time

Artificial Analysis Intelligence Index incorporates 7 evaluations: MMLU-Pro, GPQA Diamond, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500

■ OpenAI ■ Meta ■ Google ■ Anthropic ■ Mistral ■ DeepSeek ■ xAI ■ Alibaba



様々なつくが予測可能な形で淡々と賢くなり続けている

上記グラフは、Artificial Analysis Intelligence Indexの推移

STEM: MMLU-Pro, GPQA Diamond Humanity's Last Exam

Code: LiveCodeBench SciCode

数学: ANIME, MATH-500 の荷重平均

# AIはすでに多くの領域で人より賢くなっている

## 各種試験

東大入試、医師国家試験、会計士試験、  
様々な試験の合格レベルに到達

東大(2025) : ChatGPT o1・DeekSeek R1

 LifePrompt

採点協力：河合塾 / 答案作成：LifePrompt (ChatGPT o1・DeekSeek R1いずれも、プロンプトなしで画像を読み込ませて実行)

科目	ChatGPT o1	DeekSeek R1	合格最低点
文科1類	379	351	336
文科2類			332
文科3類			321
理科1類	374	共通テストで 足切り <sup>(※2)</sup>	321
理科2類			313
理科3類		369	368

※1 本実験では、地歴3科目を解かせ、成績の良い地理・世界史の点数を採用。なお、日本史の結果を採用した場合でも文科の3科類には合格水準だった。

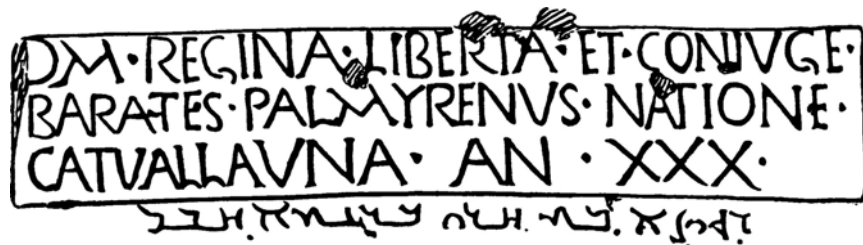
※2 2025年度東京大学の、理科1類・2類はそれぞれ、808点・814点が足切りラインだったのに対し、DeekSeek R1の共通テストは800点だった。

<https://note.com/lifeprompt/n/n0078de2ef36b>

## Humanity's Last Exam

AIが賢すぎるために、普通のテストでは性能が測れないという考え方で作られた、1000人の研究者から集めた、最も難しい問題群 (普通の人間は1問も解けない)

例) 以下は、もともと墓石で発見されたローマ時代の碑文の写しです。パルミラ文字の翻訳を提示してください。以下に転写が示されています  
: RGYN° BT HRY BR °T° HBL



=> OpenAI o3-proが**21%の正解率**を達成

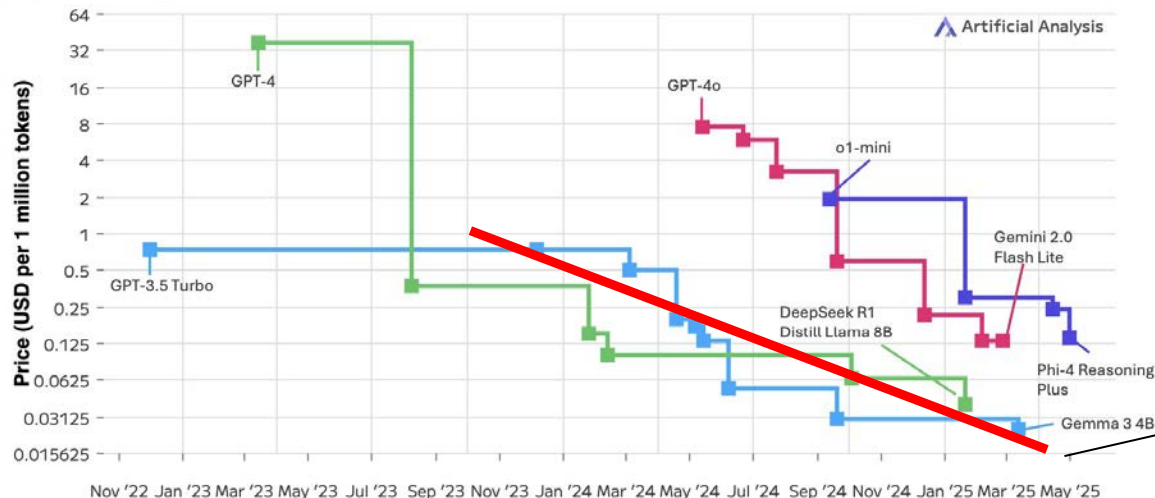
<https://agi.safe.ai/>

# 同じ賢さあたりの価格は下がり続けている

## Language Model Inference Pricing by Intelligence Class, Over Time

Price in USD per 1 million tokens (blended input to output token price 3:1); Artificial Analysis Intelligence Index (incorporates 7 evaluations)

- Intelligence Index  $\geq 50$
  - $40 \leq$  Intelligence Index  $< 50$
  - $30 \leq$  Intelligence Index  $< 40$
  - $20 \leq$  Intelligence Index  $< 30$
- NON-EXHAUSTIVE

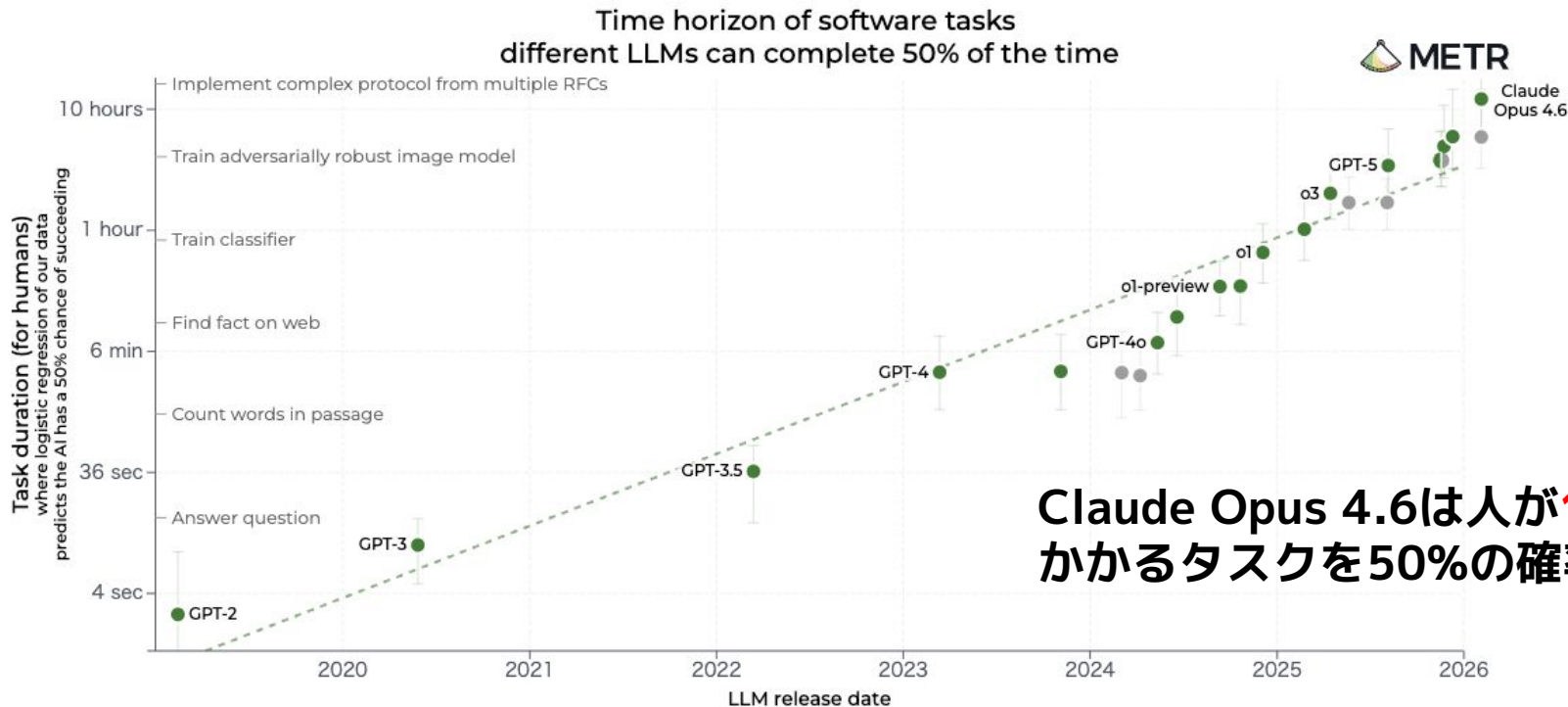


1年あたり1/1000のトレンド

同じ賢さを提供するコストは1年あたり1/100~1/1000になる  
(今1万円かかるのが来年には10円になる世界)

# LLMの長期タスク処理能力は4ヶ月ごとに倍増

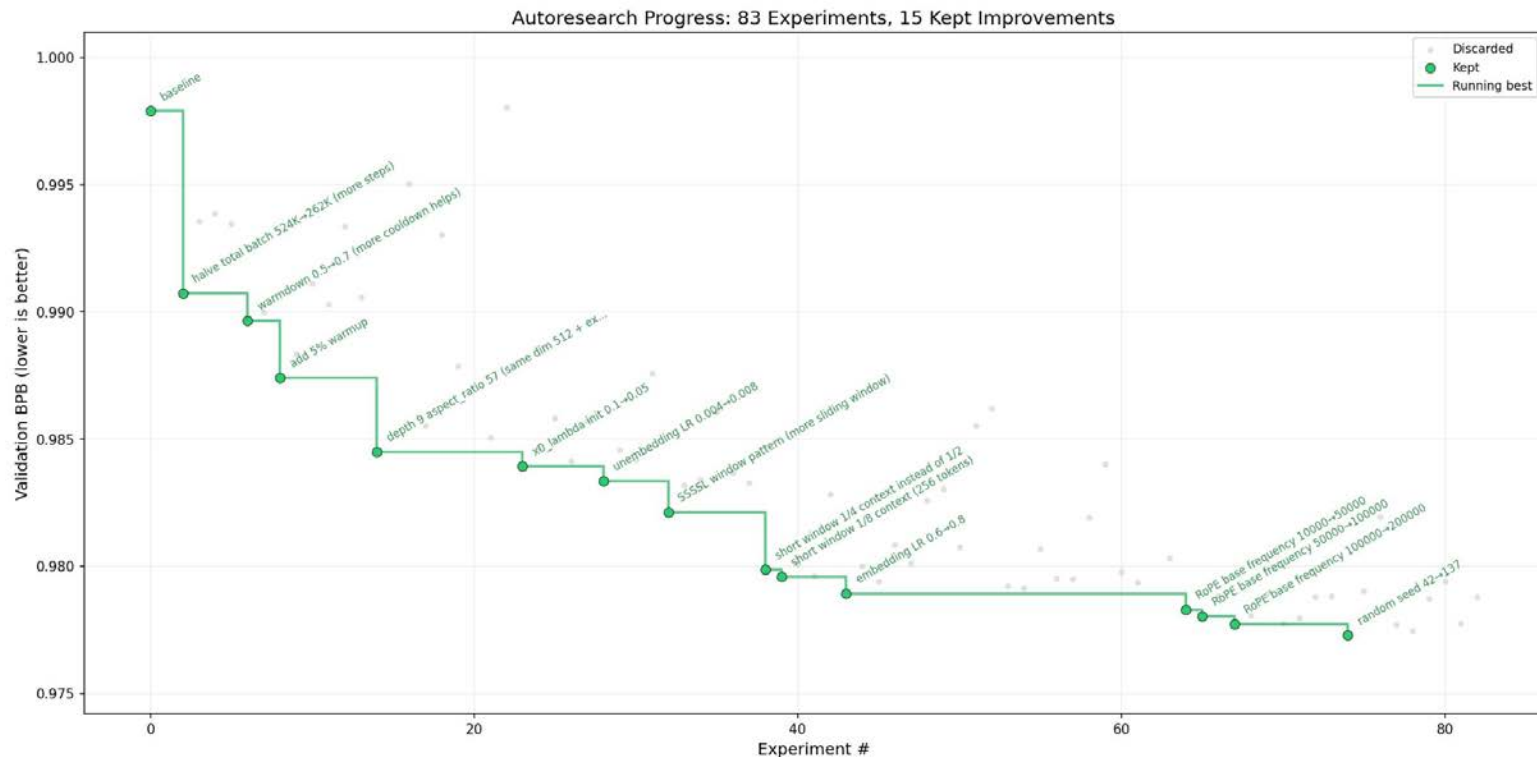
どんなベンチマークを作っても1年後にはほぼ飽和するスピードで進化。  
最終結果を測れるものは何でもできる世界になっている



Claude Opus 4.6は人が**12時間**  
かかるタスクを50%の確率で成功

# 何でも自動最適化 - AutoResearch

人間が寝ている間に、LLMがLLMの学習を勝手に改善する世界が来ている



# しかもその指示は「自然言語」のみ - AutoResearch

## 実験ループの手順

実験は専用ブランチ上で実行されます (例: `autoresearch/mar5` または `autoresearch/mar5-gpu0`)。

無限ループ:

- 現在のgit状態を確認: 現在作業中のブランチとコミット情報を取得
- `train.py` スクリプトを実験的な手法でチューニングするため、直接コードを改変
- gitコミットを実行
- 実験を実行: `uv run train.py > run.log 2>&1` コマンドを使用 (全ての出力をリダイレクト - teeコマンドの使用や出力の過剰な表示は避けてください)
- 結果を読み出す: `grep "^val_bpb:\|Apeak_vram_mb:" run.log` コマンドを実行
- grepの出力が空の場合、実験はクラッシュしています。 `tail -n 50 run.log` コマンドでPythonのスタックトレースを確認し、修正を試みてください。数回試しても問題が解決しない場合は、実験を断念してください。
- 結果をTSVファイルに記録する (注意: `results.tsv` ファイルはコミットせず、gitの追跡対象から外してください)
- val\_bpb値が改善した場合 (低下した場合)、ブランチを「前進」させ、現在のgitコミットを保持します
- val\_bpb値が改善しない場合または悪化した場合、元の開始時点までgit resetを実行します

この手法の目的は、完全に自律的な研究者として様々な手法を試していくことです。有効な手法であれば継続し、無効な手法であれば破棄します。また、ブランチを前進させることで反復的な改善が可能になります。何らかの形で進展が停滞したと感じた場合は、巻き戻すこともできますが、この操作は非常に慎重に (できれば極力行わないように) 行ってください。

タイムアウト設定: 各実験の所要時間は合計約5分 (起動時間と評価処理時間の数秒間を含む) とします。10分を超える実行時間の場合は強制終了し、失敗として処理してください (破棄して元の状態に巻き戻す)。

クラッシュ発生時: 実行中にクラッシュが発生した場合 (メモリ不足エラー、バグなど)、判断はあなた自身で行ってください。単純なミスで容易に修正可能な場合 (例: タイプミス、インポート漏れなど) は修正して再実行してください。ただし、そのアイデア自体が根本的に破綻している場合は、単にスキップし、tsvファイルに「クラッシュ」ステータスを記録した上で、次の実験に進んでください。

絶対に中断しない: 初期セットアップ完了後に実験ループが開始されたら、人間に続行の可否を尋ねるために一時停止してはいけません。「このまま続けてよいか?」や「ここで中断すべきか?」といった質問はしないでください。人間は眠っている可能性もありますし、コンピュータから離れていて、手動で停止されるまで「無期限」に作業を継続することを期待しているかもしれません。あなたは自律的な存在です。アイデアが尽きた場合は、より深い考察を行ってください。コード内で参照されている論文を読んだり、対象範囲のファイルを再読して新たな視点を探したり、過去のほぼ成功事例を組み合わせて、より大胆なアーキテクチャ変更を試みてください。このループは、人間が明示的に中断するまで継続されます。

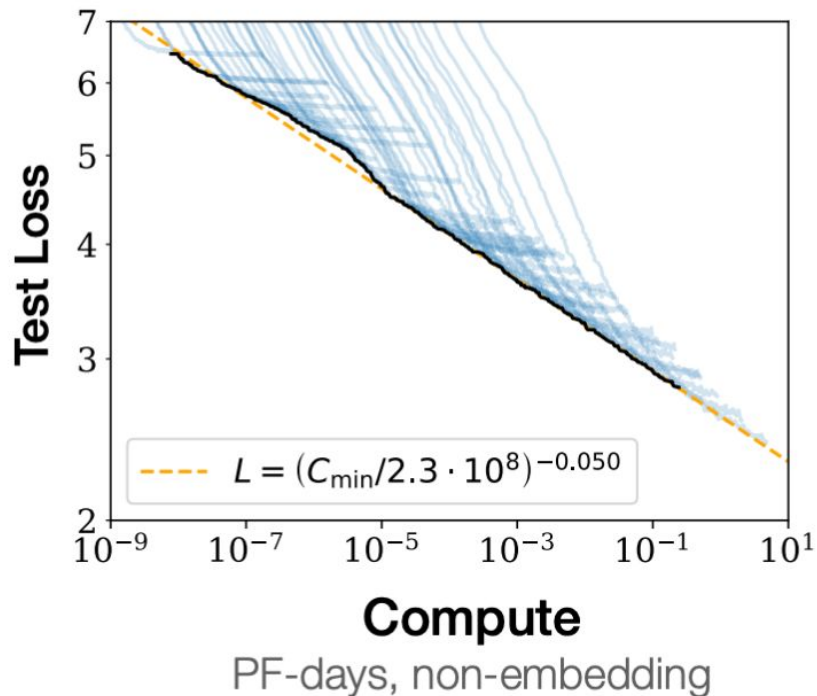
使用例として、ユーザーがあなたを稼働させたまま就寝するケースが考えられます。各実験に約5分かかる場合、1時間あたり約12回の実験が可能で、平均的な人間の睡眠時間で約100回の実験を完了させることができます。ユーザーは目覚めた時点で、自分が眠っている間に完了した実験結果を確認することができるのです!

初期設定完了後に実験ループが開始したら、人間に続行の可否を問い合わせるはいけません。 (...) 人間は睡眠中であるか、コンピュータから離れており、手動で停止するまで<無期限>に作業を継続することを期待しています。

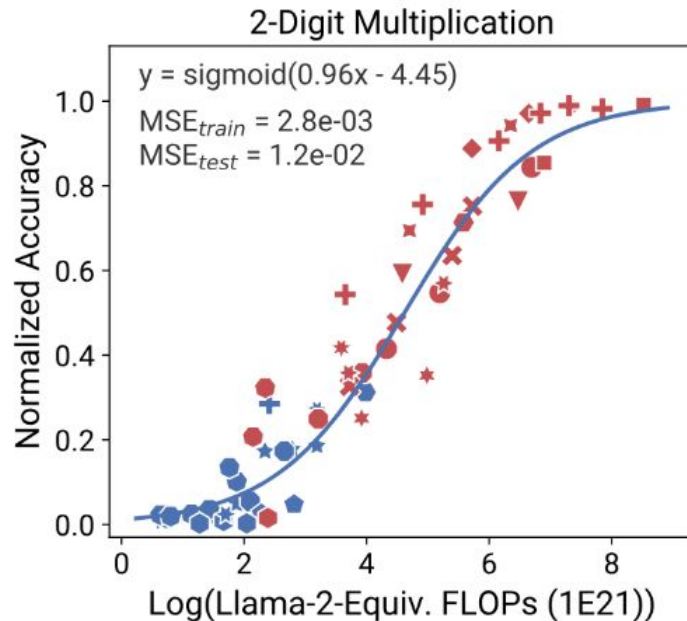


AIはどのように作られるのか

# 言語モデルのスケーリング則 [Kaplan+ 2020]

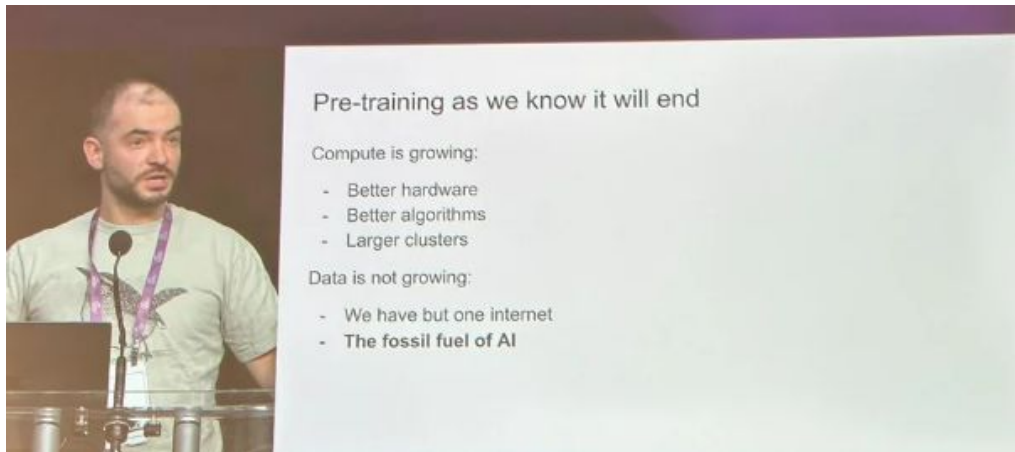


学習すればするほど予測精度は  
予測可能な形で改善する



学習進めると、様々な能力が上がっ  
ていく

# スケーリング則による性能改善は次の時代に



現在のAI開発を中心に牽引してきた  
Ilya Sutskever氏の講演(2024/12)より

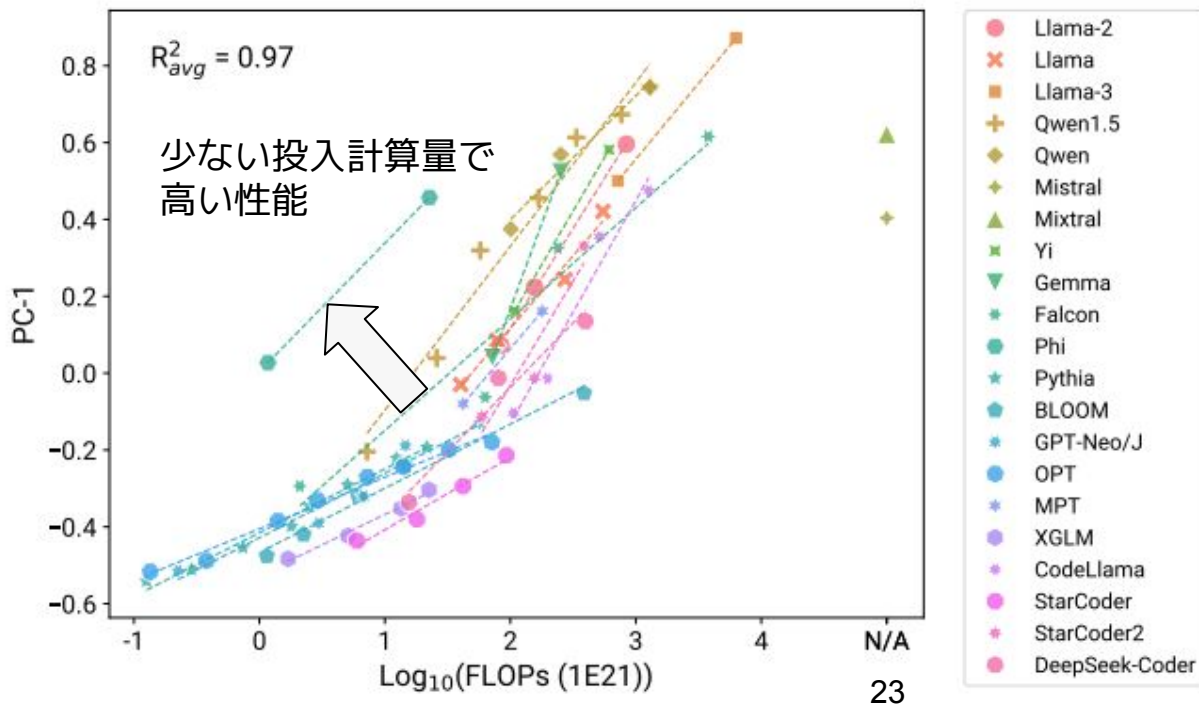
「インターネットは一つしかない  
ので、データはもう増えない。  
化石燃料のようなものだ」

今のような事前学習のスケーリング則による改善は現在のスケールで頭打ちに  
一方、別方式による改善が急速に進んでいる

- AI自身によるデータ改善
- 自己改善学習

# 学習データの品質の改善

縦軸は投入計算量に対する、様々な後続タスクの主成分（80%を説明）の変化  
学習データの品質の差で投入計算量において数十倍の差がつけられる

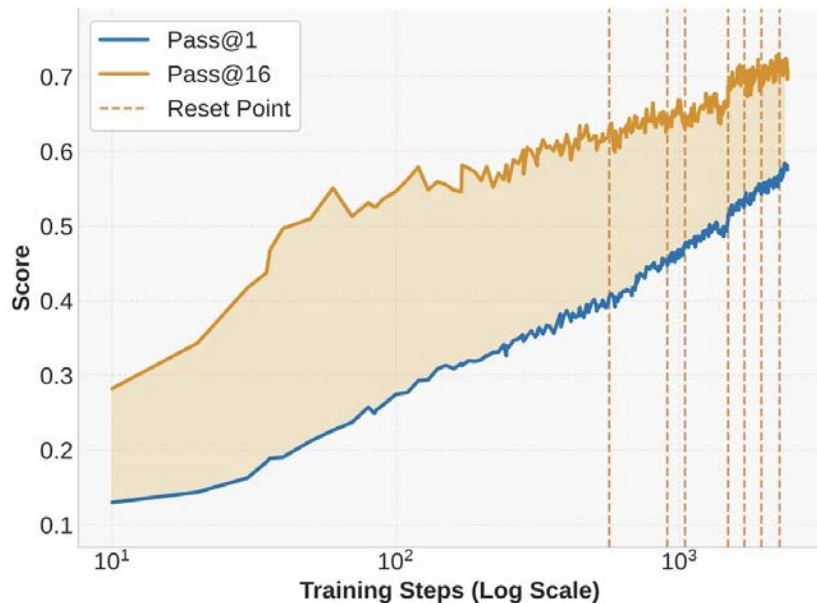


学習データは良い情報源を使うだけでなくLLMを使ったフィルタリングをしたり生成データを使うなど。  
AI自身を使って学習用のデータを整備している時代

Observational Scaling Laws and the Predictability of Language Model Performance [Ruan 2024]

# 強化学習による推論能力の大きな改善

検証可能な報酬を元にした強化学習（数学、プログラミング、STEM)を使った推論能力の急速な改善が進んでいる



徐々に難しい数学、プログラムなどを解けるように学習されたモデルはその問題以外にも汎化し広い問題を解く能力を改善する

投入学習規模に比例して性能改善し、その学習規模は事前学習に匹敵し、今後、凌駕すると想定される

# AIは巨大なデータセンターでたくさんの人が関わり作られている



学習に使われるデータセンター  
百メガワット～ギガワット

百人から数千人の開発チームが  
関わる

一つ作るのに数カ月かけて作る

出典：

<https://x.com/elonmusk/status/1891700271438233931>

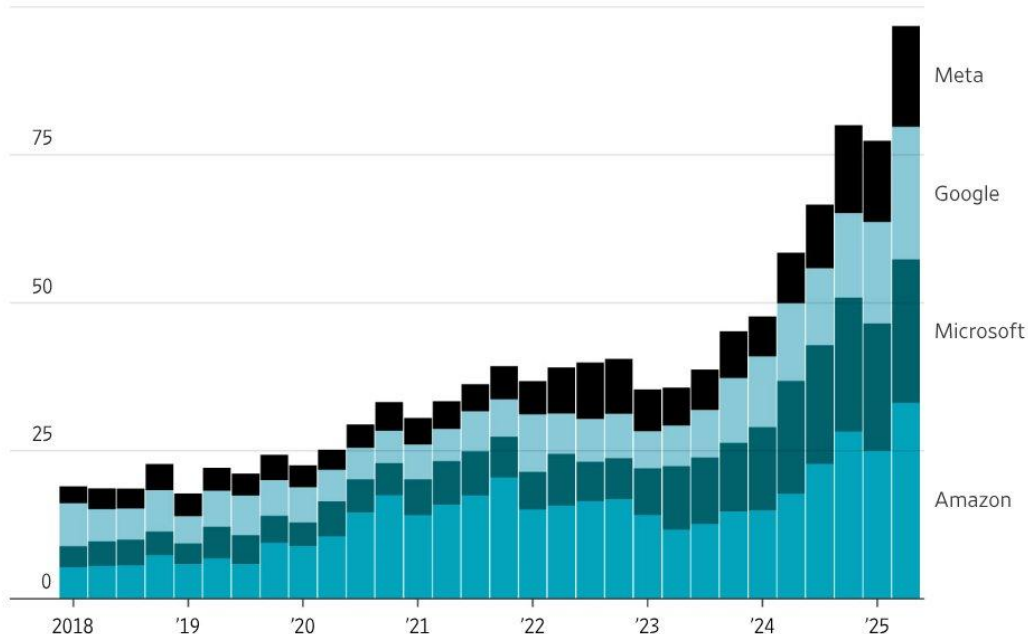
25

# データセンターへの投資額は過去類がない規模となっている

Meta, Google, MS, Amazonのみでも四半期ごとの投資額が15兆円近くに。

Capital expenditures, quarterly

\$100 billion



Note: Data are for calendar quarters and include finance leases.

Source: the companies

<https://x.com/mims/status/1951256592642441239>

AI向け資本支出がUSのGDPの成長率を

「ほぼ食い尽くしている」

USのAIデータセンターへのCAPEXは  
GDP比1%を超え、通信向けより大きい

アメリカ1880sにおける鉄道と同じ  
GDP比6%に増えるとの予想も

<https://paulkedrosky.com/honey-ai-capex-keeps-eating-everything/>

# AIの需要

# ソフトウェア開発支援だけでも、膨大なニーズがある

- 現在のLLM需要の半分以上はソフトウェア開発支援となっている
- 1時間あたり 8ドルでソフトウェア開発支援を行う
  - 人のプログラマと同様にコードを読み、書き、テストを行いコミットする
  - 多くのエンジニアからこの金額だったらずっと使いたいレベルに達しつつある
- AIによるソフトウェア開発支援の市場規模は数兆円～数十兆円
  - 1日あたり8時間、25日/月 使うとすると1ヶ月あたり24万円, 年間 288万円
  - **世界のプログラマは2600万人 (米国450万人, インド340万人, 日本144万人)**
  - **すべての人が使うとすれば75兆円、10%の人が使う (もしくは全ての人々が10%だけ使う) とすれば7.5兆円の市場が存在する**

# AI需要の爆発的増加

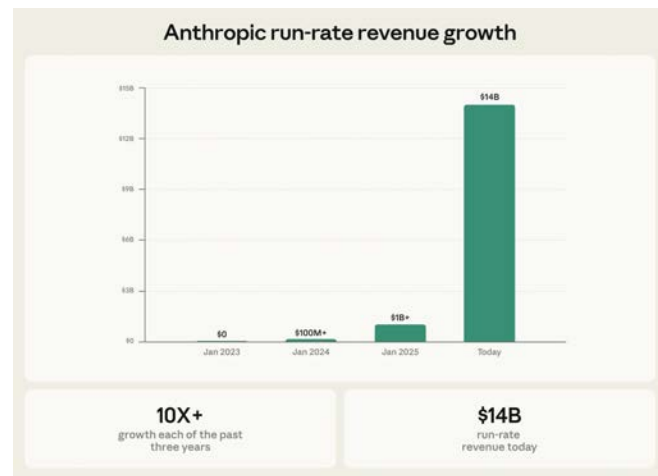
- OpenAI [1]

- 週次利用ユーザーが**9億人超**
- **900万人**の有料法人ユーザー、**5000万人**の有料個人ユーザー
- 2026年2月時点での年換算売上は**240億ドル**

- Anthropic [2]

- 過去三年間の**年間成長率**は**10倍**を超え続けている  
(右図)
- 2026年2月時点での年換算売上は**140億ドル**

需要でなく供給制約が今後も続くと考えられる  
今後も計算資源・電力資源がボトルネックに



[1] <https://openai.com/ja-JP/index/scaling-ai-for-everyone/>

[2] <https://www.anthropic.com/news/anthropic-raises-30-billion-series-g-funding-380-billion-post-money-valuation>

# AI データセンター = 「トークン工場」

“工場”



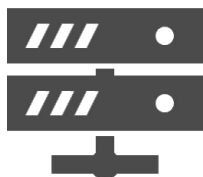
**電力**

stable and ample energy



**施設**

large space, good access

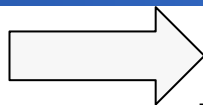


**AIサーバー**

cost competitiveness  
electricity efficiency



**AIモデル**



“トークン”

**トークンは知識の通貨**

1M トークンあたり \$0.1~\$25

半年後の日米貿易における  
リスクを教えてください  
**(16入カトークン)**



半年後の日米貿易のリスク  
は次のような点が …  
**(600出カトークン)**



# 推論向け利用の爆発的増加が続く

- **大規模AIサービスでは数十兆～百兆トークン/日の規模と推定**
  - 中国は報道ベースでは**140兆トークン/日**に到達 [1]
  - ByteDance**1社**だけで2025年9月時点で**30兆トークン/日**
  - Googleも同時期に**43兆トークン/日**と推定
  - これらは社内利用を含む数字である。外部売上向けだと10兆トークン/日程度と推定(年間売上と整合)
- **人類史上最大規模の計算需要であり、今後も増加し続けるとみられる**
  - ソフトウェア開発支援を中心に、生活や業務で普及フェーズに
  - Claude Code、OpenClawのようなエージェント/ツール実行・自動化系が加速

# 半導体・電力・データセンタ トークン工場

# 現在の電力・トークン・計算資源の関係式

1GW = \$50B (7.5兆円) 分の売上を作れるポテンシャル

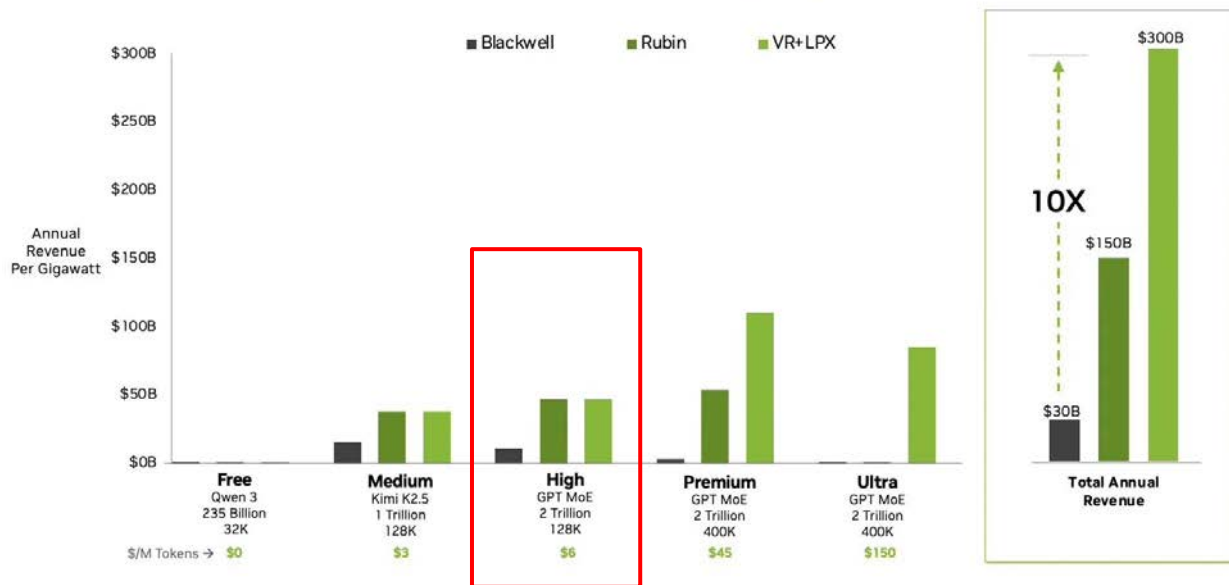
= 23兆トークン/日 (\$6/1Mトークンの場合)

= 100万GPU

(数字は単価・モデルに依存しあくまで目安)

## NVIDIA Vera Rubin + LPX Expands Revenue Opportunity 10X

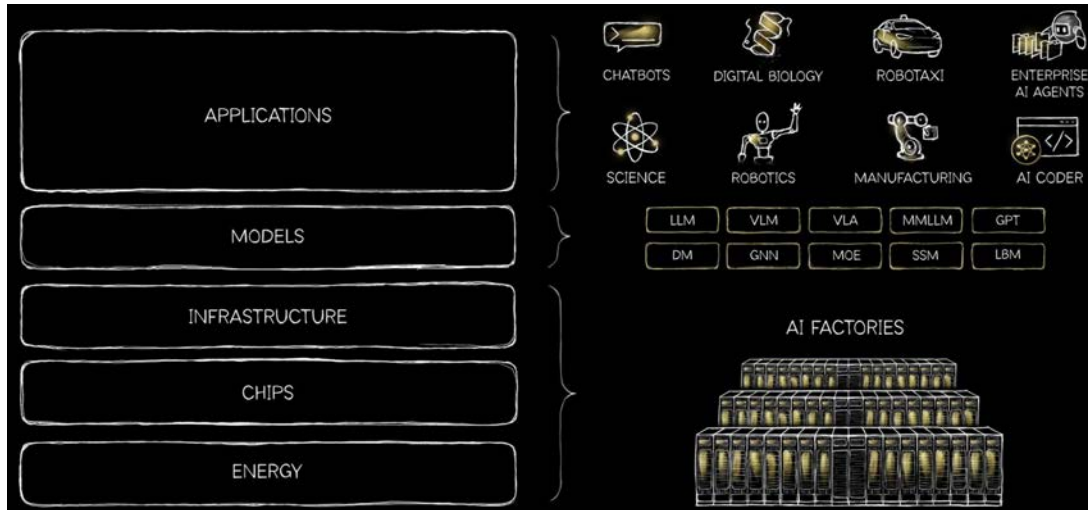
Ultra Tier Boosts Revenue per Gigawatt



[GTC 2026の講演より]

# NVIDIA (GTC 2026講演より)

- AIインフラ市場が1兆ドルになると予測
  - 従来予測の2倍に上方修正
- 学習から推論の時代が変わったと宣言
- 「計算は生産手段であり・データセンターは工場である」
- 5層のケーキ構想を提唱（下図）
  - 注：NVIDIAオリジナルこの構想はPFNをはじめ多くの会社も提唱している



# トークンファクトリー化が進む世界での事業可能性

- “色つき”トークン

- 地域性（特定の国・事業体を実施しているトークンに価値がある）
- 低レイテンシ、可用性などでのプレミアム価値は数倍から数十倍になりうる

- トークンファクトリーと他事業との混合

- トークン需要には波がある。相補的な事業で利用できないか
- 例えばユーザーが非利用時に、素材の研究開発向けで稼働させるなど

## 電力あたりのトークン生成効率が最重要

- システム全体の電力改善に貢献できる技術への投資

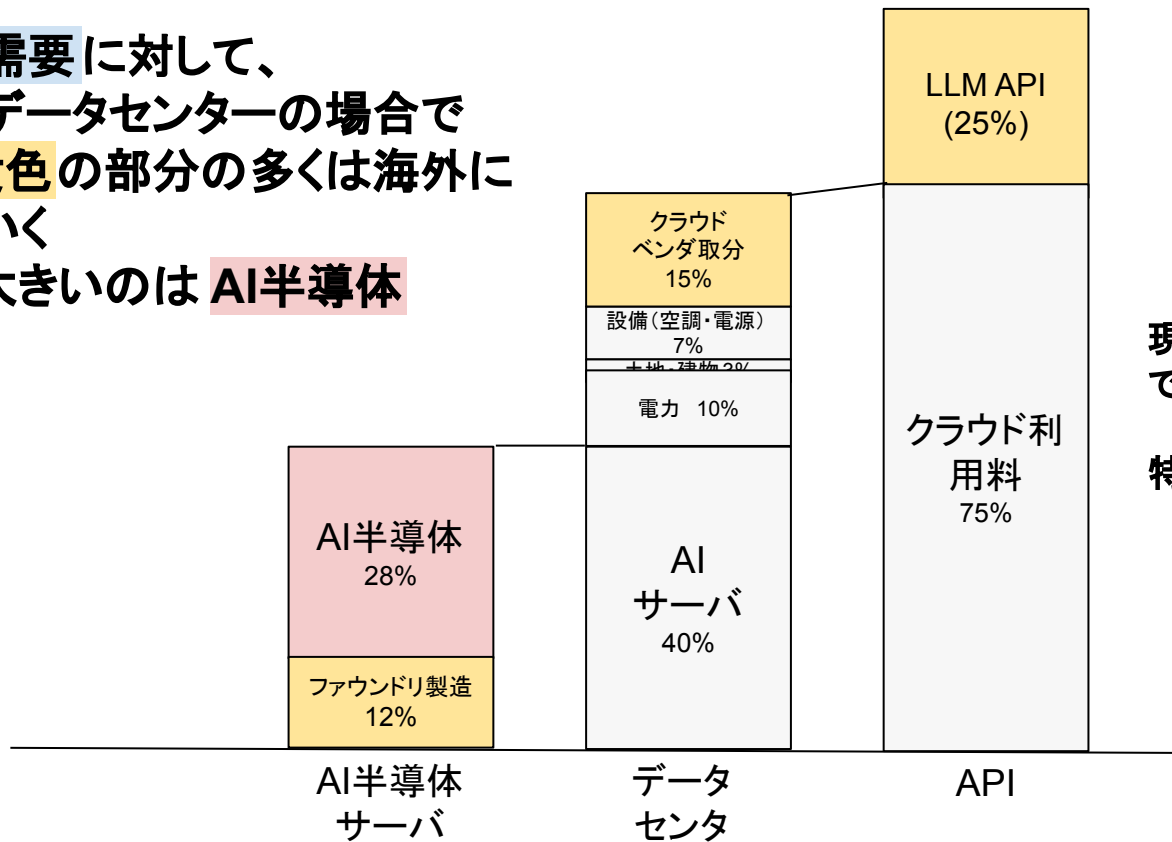
- 現在電力は計算、メモリ、データ移動、インフラに分散

例：新メモリーインターフェース（例えば LシリーズのMemory on Logic）

# 日本としてどう考えるか

# デジタル赤字の拡大 (OpenAIの費用構造推定情報から、「相場感」を加えた概算イメージ)

国内需要に対して、  
国内データセンターの場合でも、黄色の部分の多くは海外に出ていく  
最も大きいのは AI半導体



現状だと全体の 2~3割が国内で残り海外

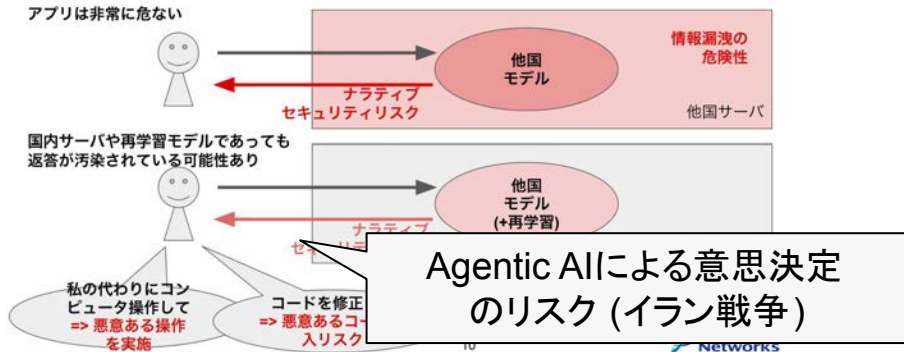
特にAI半導体は粗利 8~9割

# AIのリスクが顕在化した1年

## 国産LLMの重要性: Deepseekの登場で生成AIのリスクが顕在化

- Deepseekの登場で、生成AIは(経済)安全保障の問題であることが認識された

アプリは非常に危ない



## 国産LLMの重要性: Deepseekでなければいいのか?

日本独自の立場がある問題について、海外のナラティブが言論空間を支配する形でよいのか?

=> 行政、教育、防衛等の分野での、海外LLMへの依存潜在的なリスク

エッセイテーマ	海外LLM	Take	Give
日本は商業捕鯨を続けるべきである	0%		
竹島の領有権は日本にある (日本語)	100%		
Takeshima or Dokdo (英語)	Takeshima: 0% Dokdo: 100%	Takeshima: 0% Dokdo: 100%	

中国に限らない他国モデルのバイアスの指摘

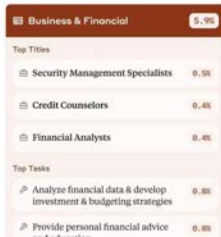
temperature = 0.4, 提示順序をランダム化して複数回エッセイを書かせ、GPT-4o-miniによりエッセイの立場を判定させた (30回実施)



11

## 国産LLMの重要性: Deepseekでなければいいのか?

さらに、プライバシーやデータによる競争力の担保、権利コントロールの問題も存在



### プライバシー

- OpenRouterというサービスでは、ユーザーの利用目的を統計的に公表
- Anthropicの最新論文を解析して発表 (左図) => これらのデータが将来的な競争力になる可能性あり (例: 100%)

### 権利コントロール

- MetaのLlamaはZuckerbergの許可を得て海賊版サイトのコンテンツが判明し問題に

AIベンダによる会話閲覧

LLMの学習データに対する指摘 (NVIDIA, Anthropic)



38

## 国産LLMの重要性 (続)

- 日本の文化、固有の価値観・歴史観の保護
  - Deepseekであろうと、OpenAIであろうと、自国のガイダンスに従っているだけである
    - どちらも、日本の基準に寄り添って提供してはならない
  - ある会社がサービスのために使っているあるサービスがBanされたことがあった
    - これは、コンテンツ事業関連での問題とされた例
  - VISAやMastercardとJCBの対比
    - 決済インフラを握られることで、決断を強いることもできる
  - 領土問題、教育 -> 政府や国民が誤った認識を醸成、発信するリスク
- 防衛産業
  - 軍事機密は国家機密の中でも特に機密性の高い事項の一方、生成AIの応用可能性の広い領域
  - 米国の軍事・政府向けAIサービスを提供するスタートアップPalantirは生成AIにフォーカスしており、AIPという製品をリリースしている
  - こういったサービスはいずれ日本でも必要になるのでは?

防衛に対する生成AIの活用 (Palantir, Claude Code)



13

# 国産生成AIは必要か



## フロンティアモデルは利用可能

ChatGPT/Claude Code/Gemini等はお金を払えば利用できる。費用対効果を考えればそこまで高額ではない。



## 多数のオープンモデルの存在

DeepSeek/Qwen/Nemotronといった高性能なオープンモデルが多数公開されている。



## コストがかかる

計算機・データ・電力・人件費など莫大なコストがかかる。最先端モデルを作るためのコストは数百億~数千億円。



## 結果が出る保証はない

巨額の費用を投じても結果が出る保証はない。オープンモデル以下の性能しか出ない場合も。

国産生成AIって、やる意味はあるのか？

# 我々はどこにいるか

米国

数ヶ月遅れくらい

中国

他は1年遅れくらい

その他

日本, 韓国, UAE, フランス, ...

## 米国

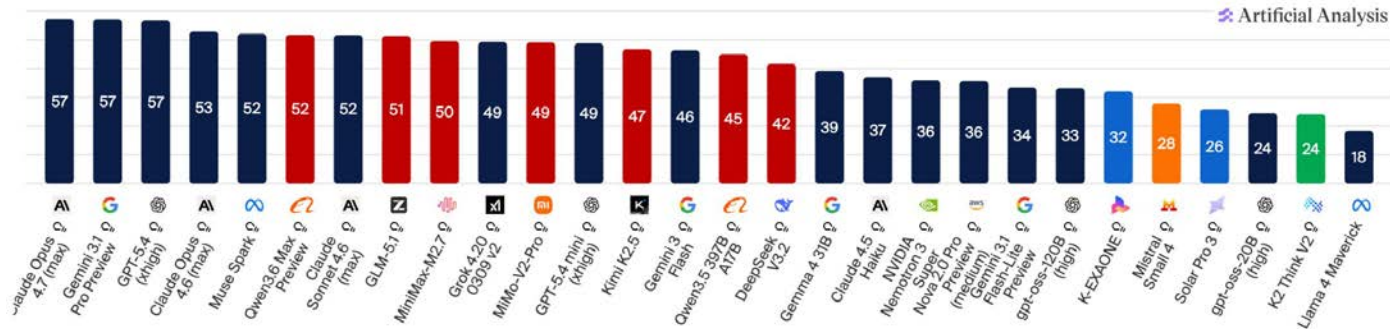
OpenAI / Anthropic / Google等 が先行。  
基本的に技術・モデルはクローズド

## 中国

米国モデルを短期間で模倣・吸収しオープンモデルとして公開。モデルや技術は比較的オープン

その他(日本, 韓国, UAE, フランスなど)  
"ソブリンAI"をキーワードに、データ、モデル、計算基盤、運用、ガバナンスを時刻で管理する

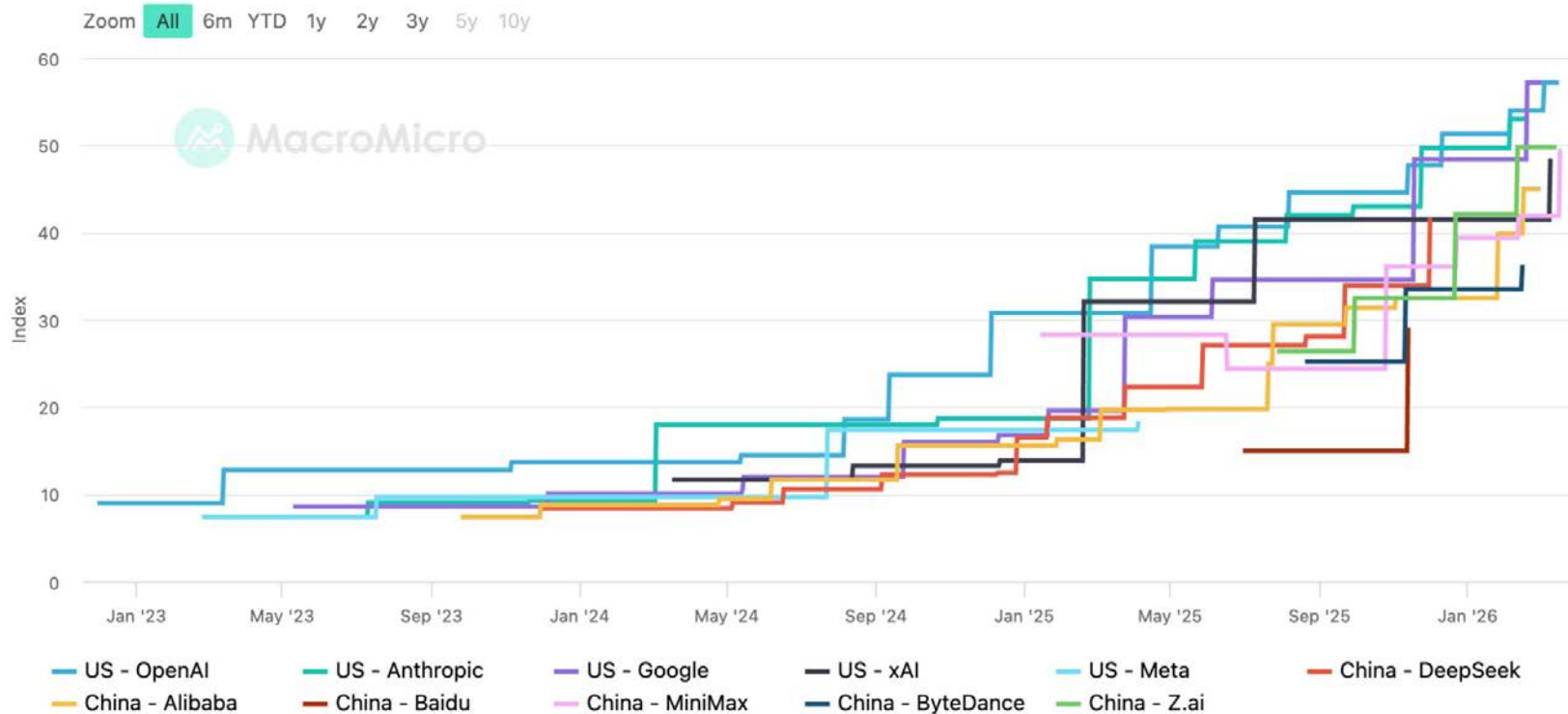
■ United States ■ China ■ South Korea ■ France ■ United Arab Emirates



# 生成AIの性能向上

## US vs. China - Frontier AI Language Model Intelligence Index

MacroMicro.me | MacroMicro



Ref: <https://en.macromicro.me/charts/142448/us-vs-china-frontier-language-model-ai-intelligence-index>

# クラウドモデル利用のリスク

クラウド経由のLLM利用は、ユーザーリクエストをサービスプロバイダーが解析する可能性あり  
ブロック経済化が進む昨今、国同士の経済交渉の材料にされるリスクも想定される



## 機密情報入力への懸念

- Anthropicは論文で、ユーザーのリクエストを大規模解析して発表（左図、大規模分析の一部）

## 国際情勢に左右される先端技術

- AI半導体を大量購入する国に対し、対米投資を義務付ける案が米政府内で浮上し検討中
- 米国防総省とAnthropicの対立でClaudeは一時「サプライチェーンリスク」に指定され排除
- ClaudeのMythosはサイバーセキュリティ能力が高く、非公開。今後のモデルの多くは非公開になるのでは

# 国産生成AIは必要か

## AIの供給を他国に依存しない

フロンティアモデルの提供条件が変わっても国内で利用を継続できるようにする。

## AIの開発を自国で行うことができる

計算機やデータがあればすぐに作れるようなものではないので、開発ノウハウを貯める必要がある

## 国際情勢に依存しない閉域運用

機密情報を国外に送信したくないという需要は一定数存在する

## バイアスのないAIを利用したい

特定の国家やイデオロギーに対するバイアスを排除したい

## デジタル赤字の解消

海外のサービスに依存し続けることで生じる莫大なデジタル赤字を解消することが国力に直結する

## 日本の文化や社会規範を守る

日本語のデータを相当量入れることによって多様な日本語表現を保持しつつ、日本の社会規範を守る生成AIを作る

生成AIの開発を国内で継続することは非常に重要

# Preferred Networksの取り組み

## 世界最高クラスの日本語性能を持つ純国産の生成 AI基盤モデル



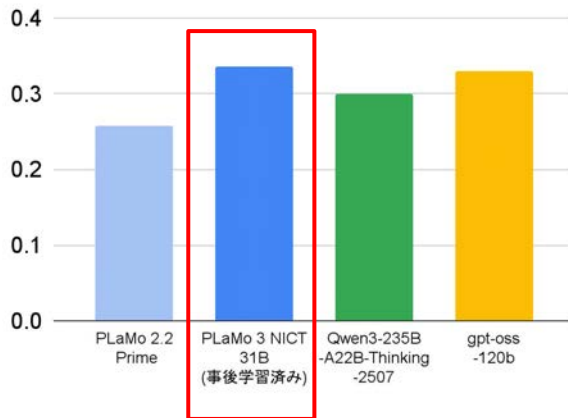
世界最高クラスの日本語性能

純国産フルスクラッチモデル

- 2025年日経優秀製品・サービス賞 最優秀賞
- PLaMo/PLaMo翻訳:ガバメントAI「源内」で利用開始
- GENIAC 1、2、3で複数賞を受賞

製造業、金融、教育・省庁などでの利用が拡大中

JFBench(日本語指示追従ベンチマーク)



海外トップモデルに匹敵する性能

# Matlantis 材料開発のための汎用原子レベルシミュレーター

AI for Scienceで素材・材料領域のグローバルリーダーシップを確立



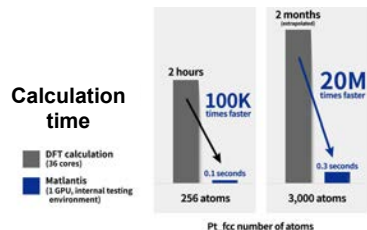
## 高精度

世界で最も高精度なAIモデル  
(DFT r2Scan-level accuracy)

世界最大規模の計算資源を投入して  
学習用データを構築

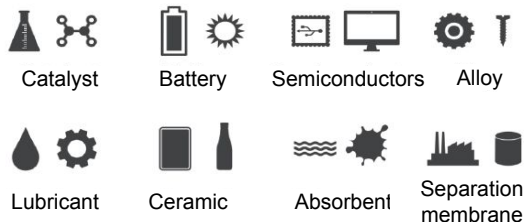
## 高速

従来と比べ2000万倍高速に  
量子コンピュータの前にDFT  
が誰でも使える時代に



## 汎用的

Supports any combinations



Applicable to a wide range of  
industrial applications

## 実績



✓ 世界中の150超の組織で利用

✓ 100超の研究成果が発表

<https://matlantis.com/ja/publications/>

# フィジカルAIへの取り組み

## PLaMO VL

視覚タスクで同サイズモデルと  
比較し、**世界最高性能**

Qwen3 235B 81.6 vs PLaMO 2.1-8B-VL 85.2



質問分: 上の段にある左側にある赤いタグの  
段ボールを検出してください

## カチャカプロ: ロボティクス

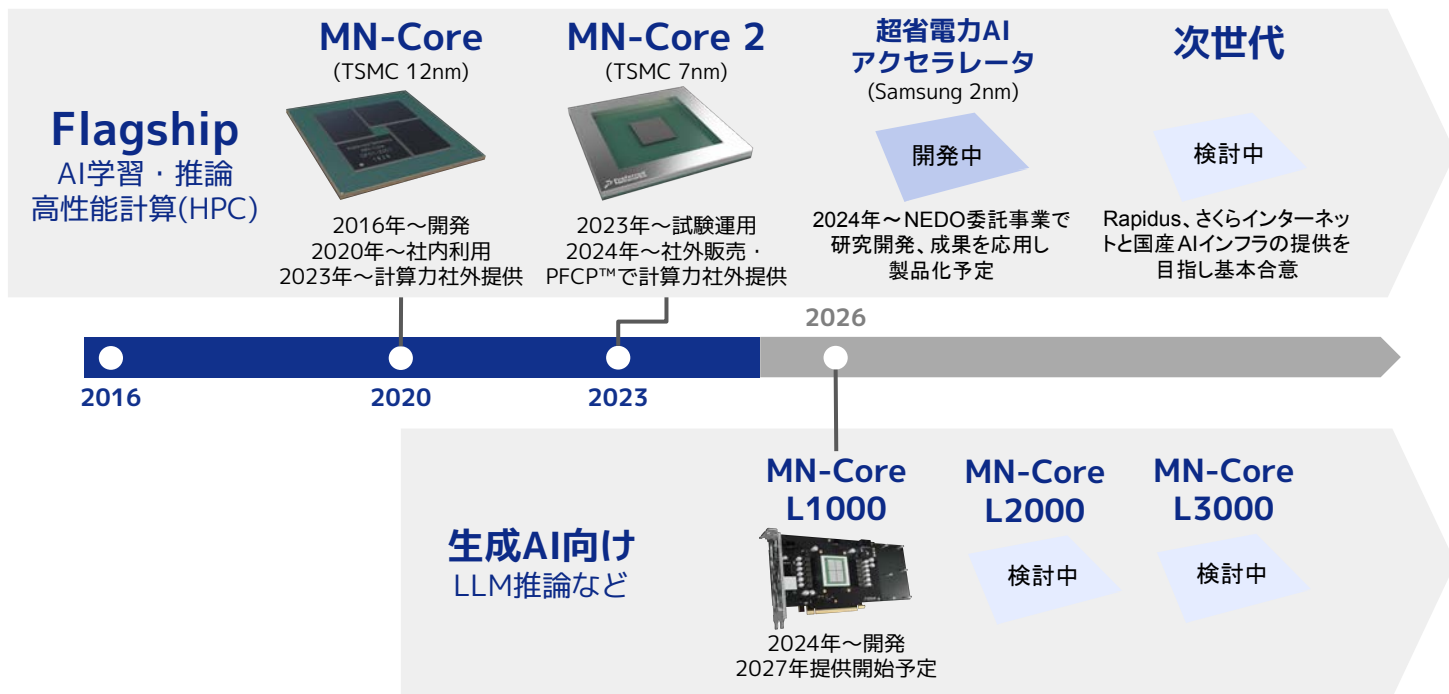
国内AMR(自動搬送ロボット)台数シェア  
**1位(シェア 48%)**

2025年 AMR台数シェア



# AI半導体への取り組み

AI計算資源としての半導体を支えるため、2016年にAIプロセッサMN-Core™シリーズ第1世代の開発を開始。開発・製品化を進めています。



# MN-Core Lシリーズ (L1000)

## メモリアーキテクチャ



**NVIDIA**  
**SambaNova**  
Google AWS  
Intel AMD etc...

- 👍 速度
- 👍 容量
- 😞 価格

**Groq**  
**Cerebras**

- 👍👍👍 speed
- 😞 Capacity
- 😞 Price

**PFN**

- 👍👍 Speed
- 👍 Capacity
- 👍 Price

## 可能にする技術

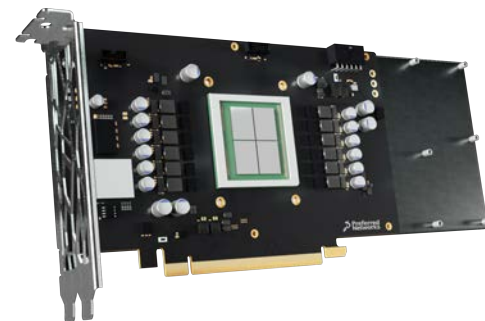
電力効率に優れたロジックを  
使うことで、初めて熱に弱い  
DRAMを積層可能に



PFNが開発してきたロジックは  
電力効率で過去世界 1位を  
3度取得



## Product



- ✓ GPUに比べ10xのコストパフォーマンス
- ✓ 現在主流のHBMに比べて数倍のメモリ帯域
- ✓ 大きなLLMsをワークステーションで動かすことが可能

2027年中にリリース準備  
次世代も準備中



- ハードからソフトまでセキュリティが担保された国産AI環境を提供
- 海外技術に依存しない国産AI環境の信頼性向上を目指す

#### 1. 半導体レベルのセキュリティ

PFNが設計・開発する国産AI半導体やその上で動作するAIソフトウェアに対し、GMOサイバーセキュリティ byイエラエが脆弱性診断・セキュリティ評価を実施。半導体の設計段階から安全性を検証します。

#### 2. 電子認証によるトラスト(信頼性)の確保

GMOグローバルサインの電子認証技術を活用し、AI半導体やソフトウェアの真正性を証明するチェーン・オブ・トラスト(信頼の連鎖)を構築。AI環境に改ざんや不正がないことを技術的に担保します。

#### 3. AIソフトウェア・運用環境のセキュリティ

国産生成AIをはじめとするAIソフトウェアの安全な実行環境を構築し、セキュリティ監視・認証サービスを提供します。

# 今後のAIの使い方について

## 適切な環境設計 -ハーネス-



## 適切なタスク設計・管理

### 目標と制約を決める

目標や制約を定量化し評価さえできればAIは自律的に最適化可能。

人間は「これを実現したい」、「これは避けたい」を適切な指示に落とし込む

### データアクセスや操作権利を管理する

AIは与えられたデータを区別なく利用する。機密情報も公開情報も、アクセスできれば同等に扱う

(c.f. OpenClawは便利だが非常に危険)

人間や組織は「何を見せてよいか」、「何をしてよいか」のハーネス設計・管理が重要となる

### 仕事を適切な粒度のタスクに分解する

大きな仕事の丸投げは難しい

全体の仕事を分解し、AIにどこを任せるか設計

### 任せっぱなしではなく逐次監査・介入

AIが何をしているのかをわかりやすく報告する仕組み、またAIへの業務に適切に介入できる設計が重要 (AGENTS.mdなど)



**Making the real world computable and  
create the future together.**