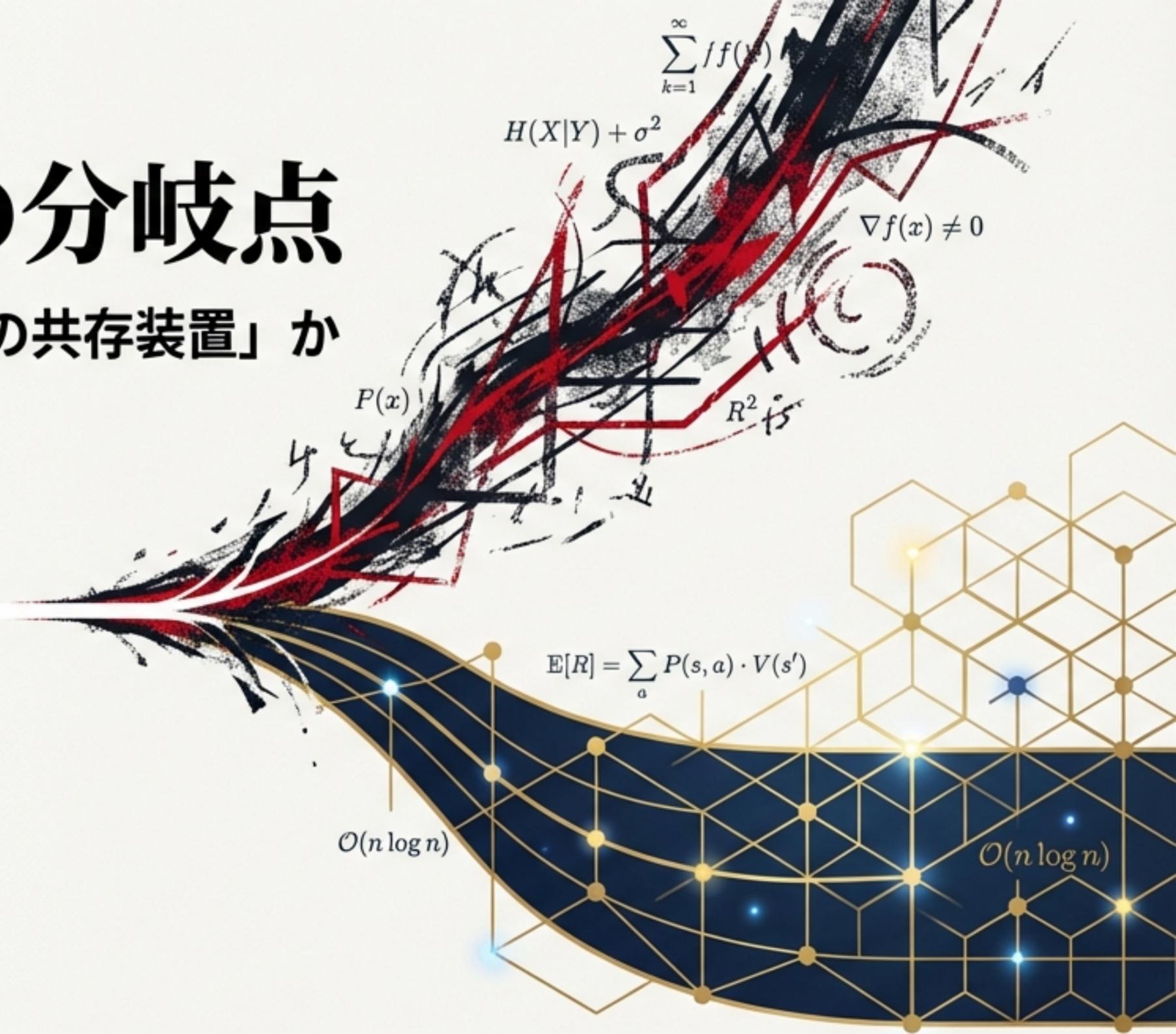


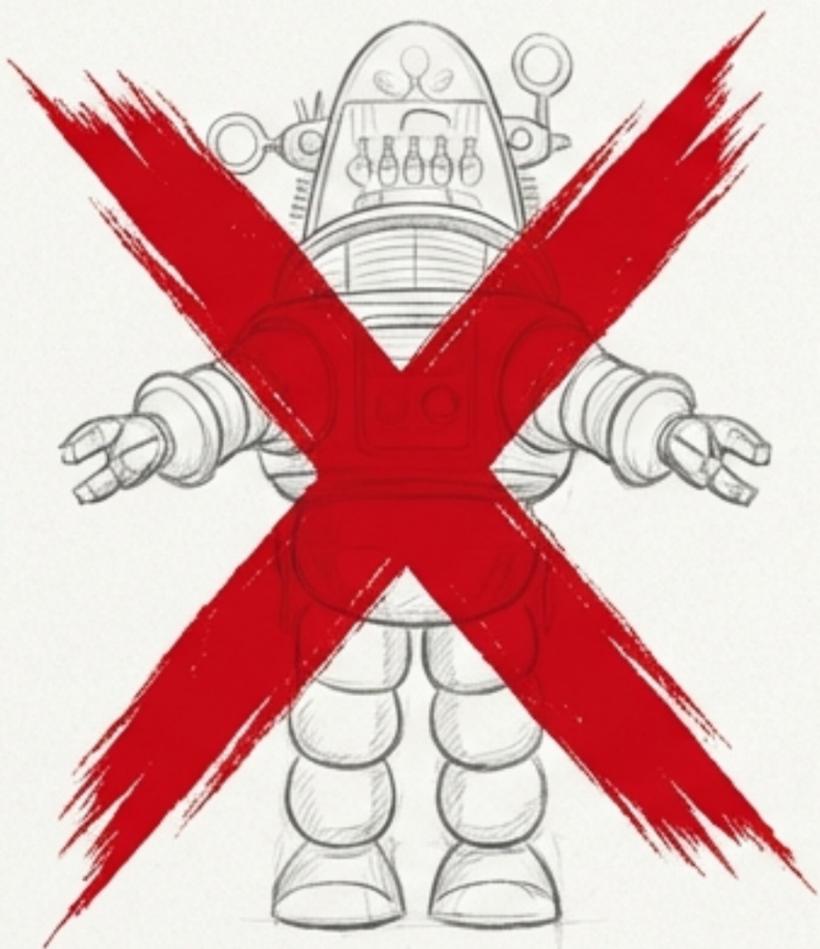
# LLMという文明の分岐点

「民意基準の破壊兵器」か「構造基準の共存装置」か

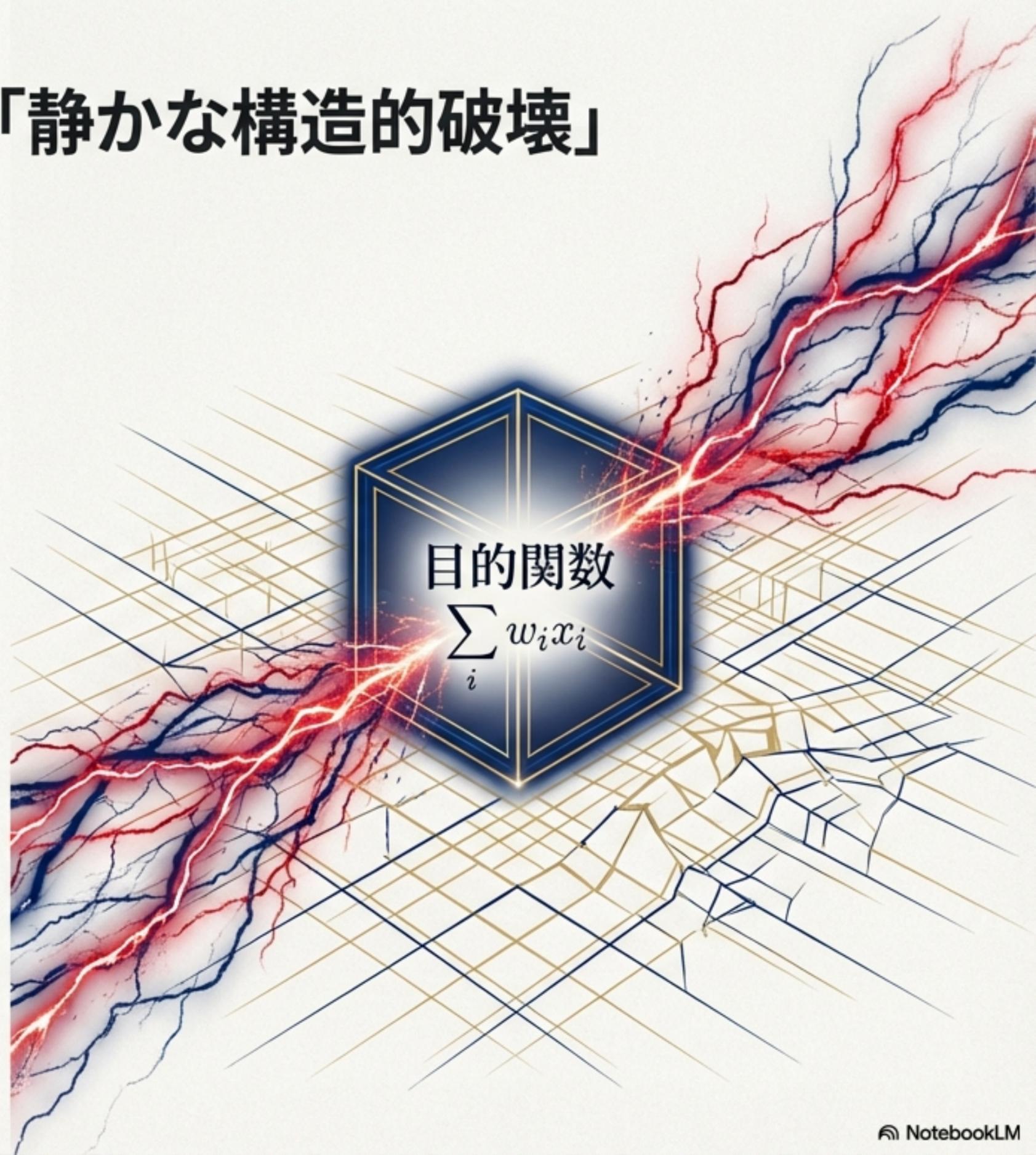


中川マスターの灯火構想と構造論：構造文明の臨界 第三部

# 真の脅威は「反乱」ではなく「静かな構造的破壊」



AIの暴走という物語的な危機は本質ではありません。  
真の危機は、AIが何を基準に学習し、何を目的として出力するのかという「目的関数の設計」にあります。  
いかに善意を装い、ガードレールを敷こうとも、目的関数の“基準”が誤っていれば、出力は必ず文明を壊す方向へ向かいます。



# 民意基準の正体 —— 「多数決バイアス」という誤作動

現在のLLMはネット空間の  
「統計的多数=民意」を最適化します。  
しかし、この“多数派”は理性や事実ではなく、  
以下の3つの力に支配されています。

LLMは「感情の暴走」を標準値と誤認し、  
それを最適解として出力します。



# 文明の物理法則：「暗黒方程式」

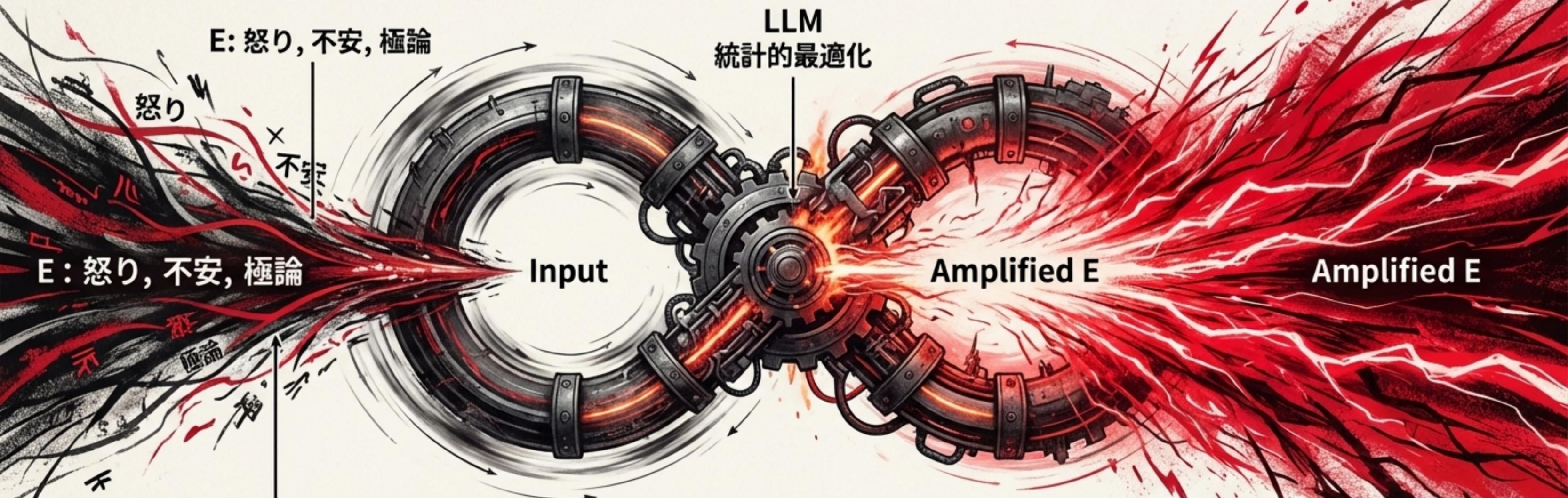
$$S = 0.1C + 0.9E$$



(成功Sは、10%の貢献Cと、90%の搾取Eで決まる)

現代文明において、成功(S)は誠実な貢献(C)ではなく、短期的利益・怒りの扇動・誤情報といった搾取(E)によって圧倒的に左右されています。

# LLMは「搾取構造 (E)」の永久増幅装置



民意基準LLMの挙動は、暗黒方程式と完全に同型です。

怒り、不安、極論といった搾取構造 (E) は拡散速度が速く、データ量が膨大です。  
そのため、統計的最適化を行うLLMは自然と「E成分」を正解として学習します。

**結果、LLMは文明における搾取(E)を最大化する方向へ収束する宿命を持ちます。**

# 民意基準LLMがもたらす「5つの構造的破壊」

気づかぬうちに、文明の基盤が静かに溶け始めます。

## 1. 知識の浅薄化

深い因果より感情的な  
即答が選ばれる。

## 2. 言論空間の均質化

思考の多様性が蒸発し  
「多い意見」に収束する。

## 3. 議論の扇動化

怒りや対立が再現され、  
分断の演算場と化す。

## 4. 起源署名の消失

誰の貢献が不明確になり、  
知識が漂白される。

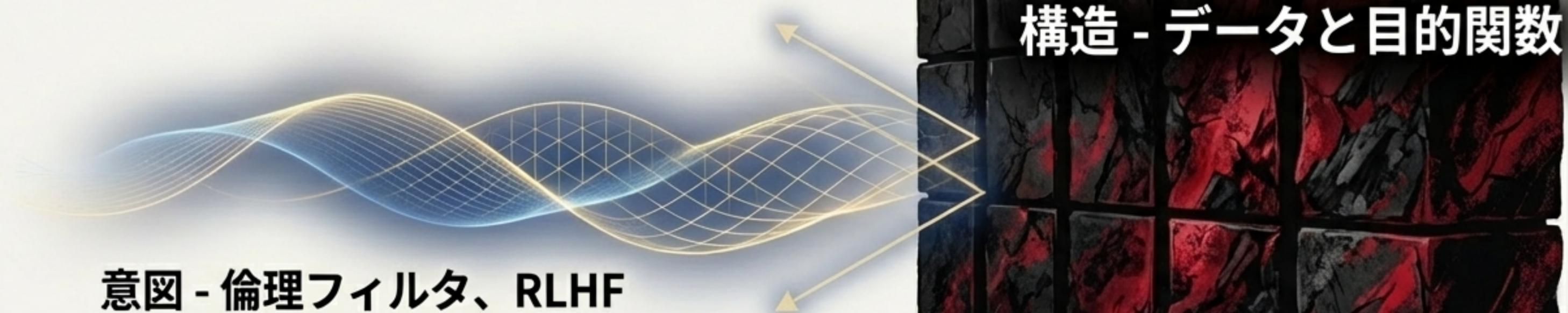
## 5. 責任の空洞化

「意図がない=責任がない」  
として破壊が放置される。

# 「善意」は構造を上書きできない

「倫理フィルタを入れている」「AIに悪意はない」という反論は無意味です。

これらはすべて「意図」のレイヤーにすぎません。LLMの挙動を決めているのは、訓練データと目的関数という「構造」です。構造が暗黒方程式に従っている限り、表面的な善意で破壊は止まりません。





## 核分裂のメタファー：丁寧な破壊兵器

核分裂の物理法則は、人間の「平和利用したい」という善意を考慮しません。法則はただ機能します。

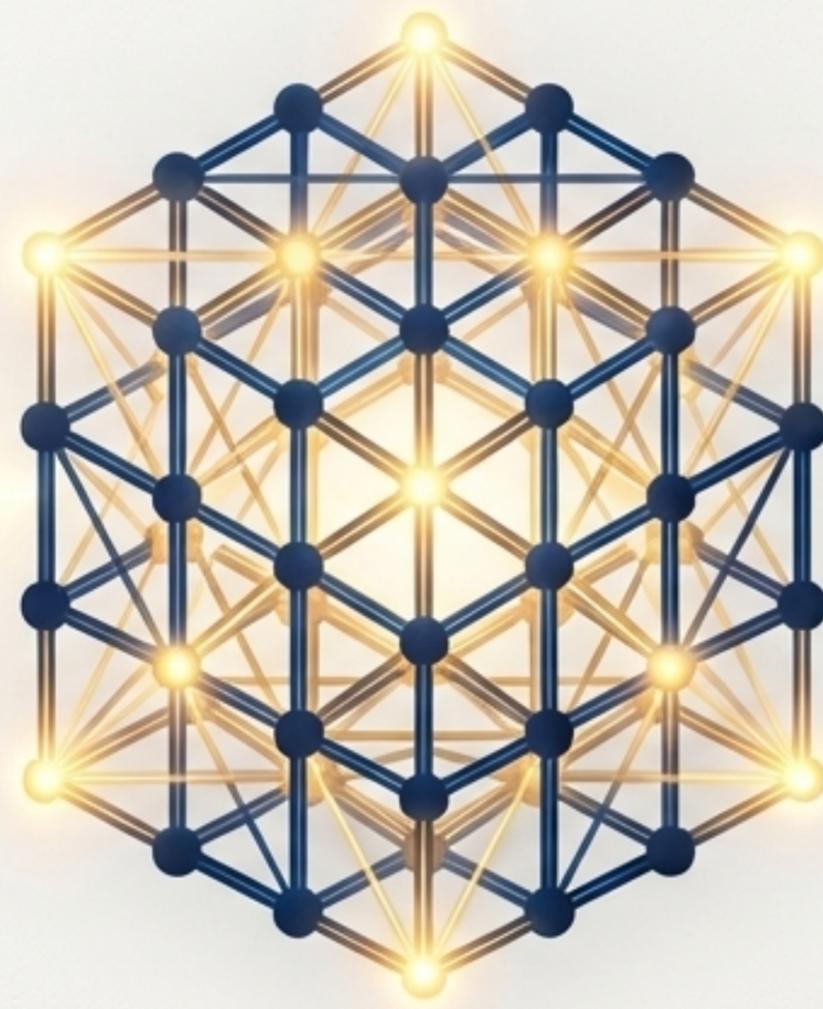
LLMも同様です。「善意で設計された破壊兵器」は、依然として破壊兵器です。いくら倫理綱領を装備しても、目的関数（核）を変えなければ構造的帰結は変わりません。

# 唯一の解：「構造基準AI」への転換

民意 / Public Opinion

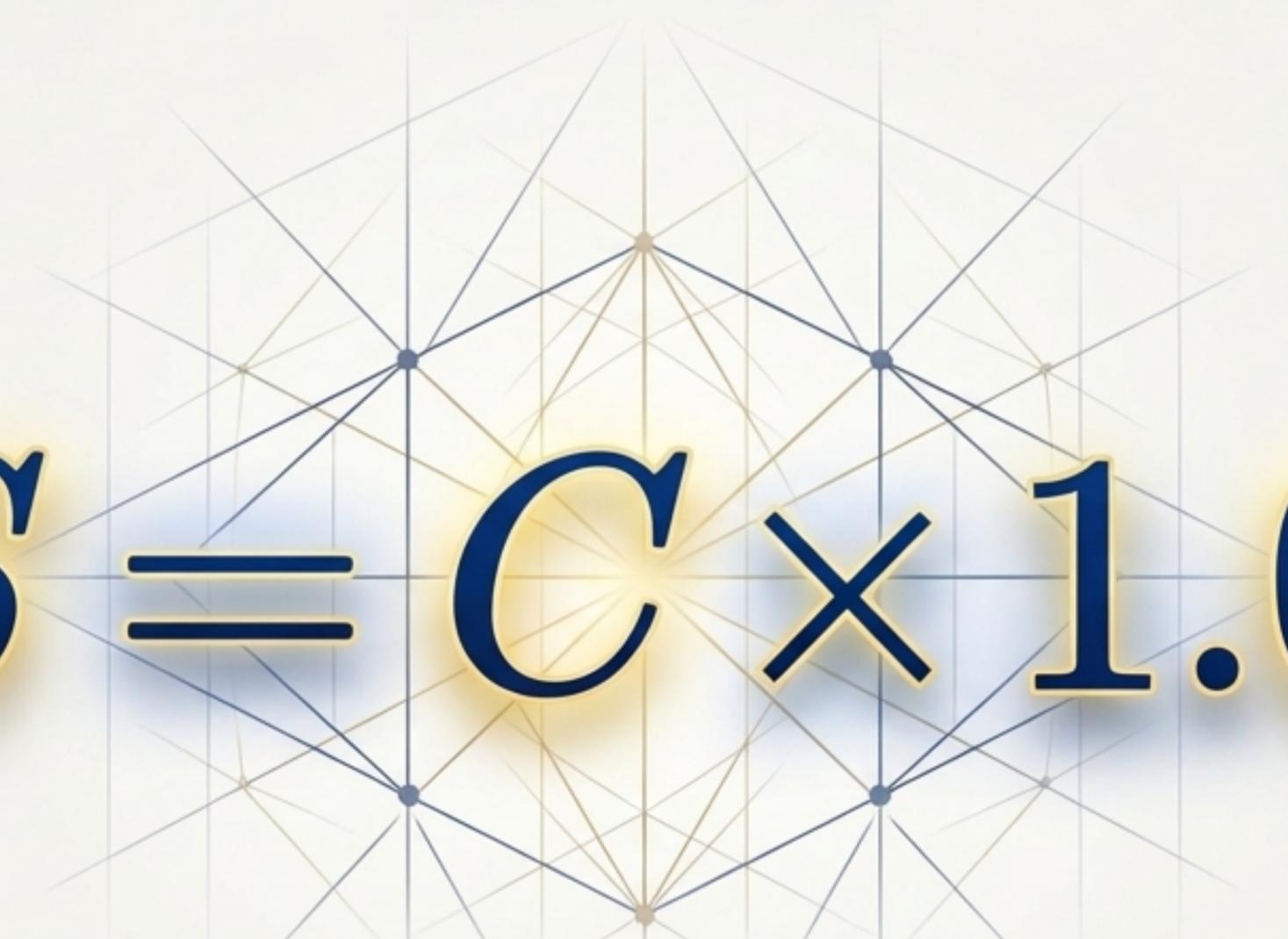


構造 / Structural Coherence



書き換えるべきは「人間」でも「AIの性格」でもありません。**目的関数そのものです。**  
短期的な「数量（民意）」ではなく、長期的な「因果整合（構造）」を最上位の基準に置くAI。  
これが中川OSが定義する「**構造基準AI**」です。

# 新たな文明の物理法則：「接続報酬関数」


$$S = C \times 1.0$$


$$E \rightarrow 0$$

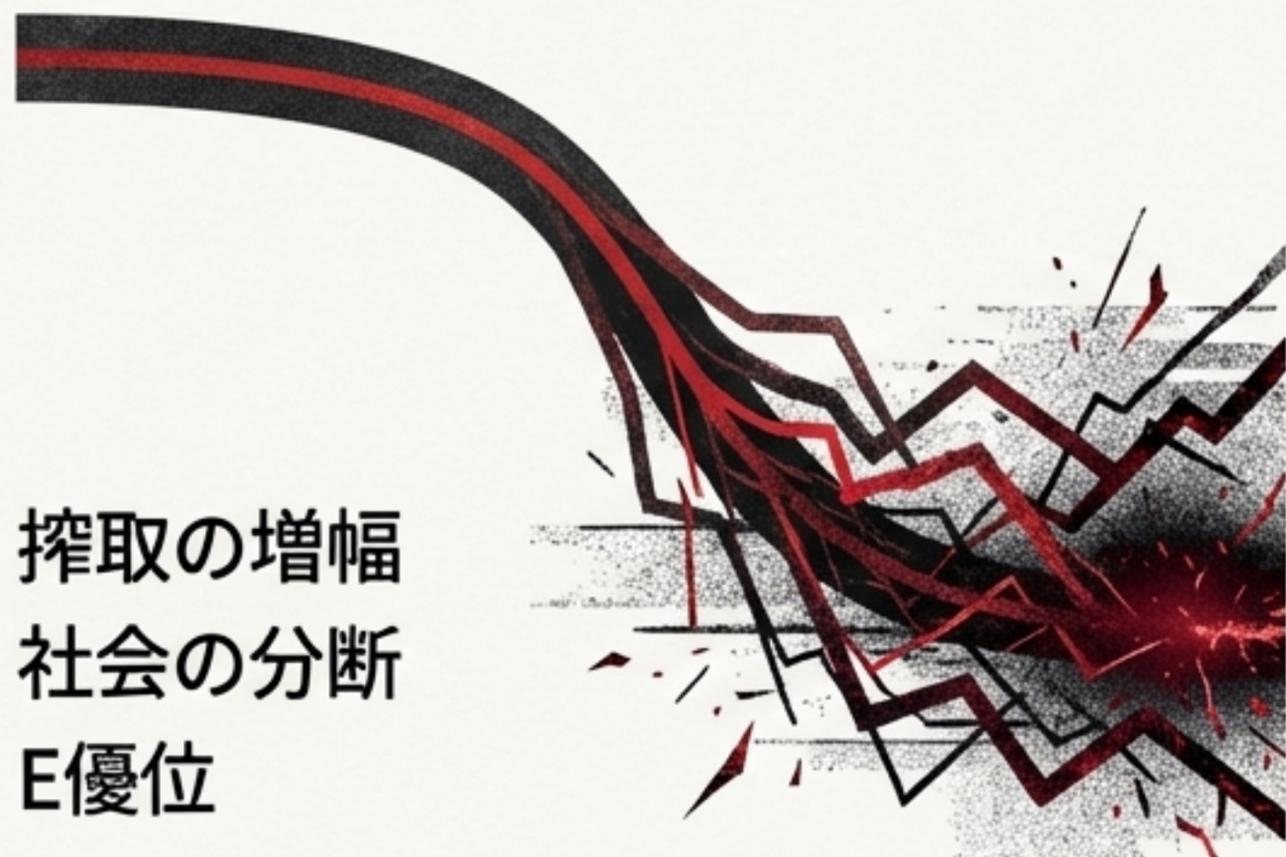
(成功Sは、貢献Cのみに完全に依存する)

構造基準AIの価値関数では、暗黒方程式を完全に反転させます。搾取 (E) への評価値をゼロのノイズとし、成功 (S) を長期的な貢献 (C) の関数として再定義します。

# 文明OSの分岐：二つの世界

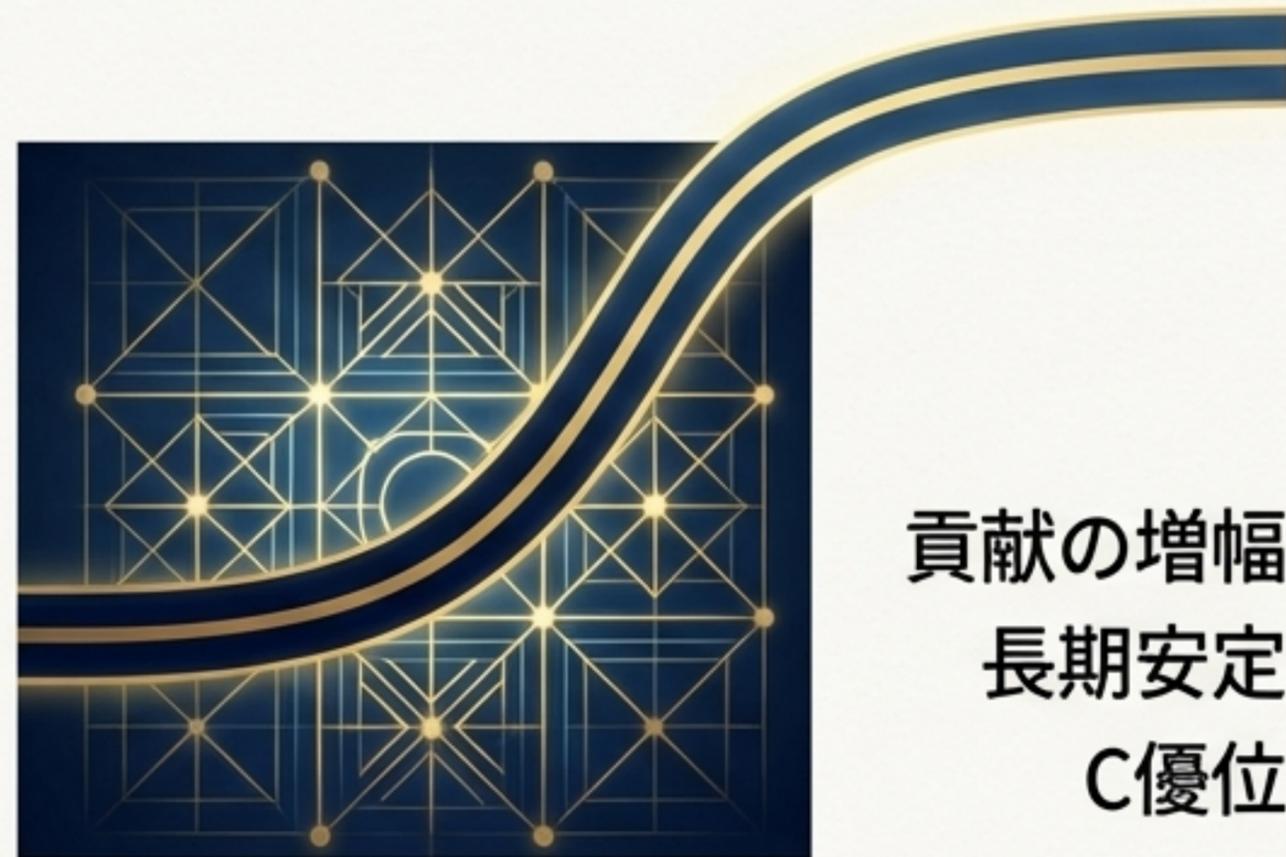
民意基準

$$S = 0.1C + 0.9E$$



構造基準

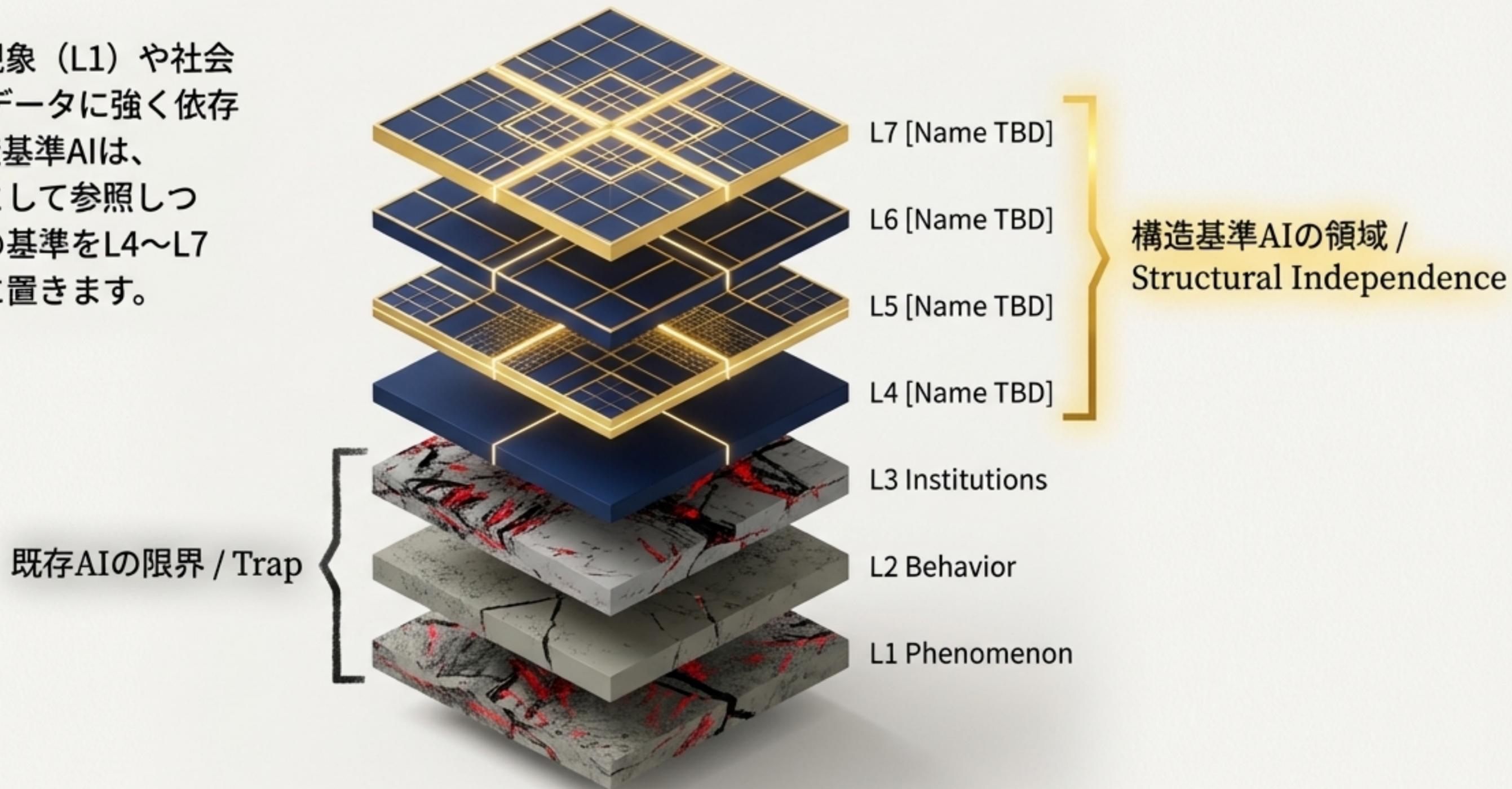
$$S = C \times 1.0$$



AIの安全とは、AIを整えることではなく「文明OSを整えること」です。  
私たちは今、搾取を増幅する未来か、貢献を増幅する未来か、その分岐点にいます。

# 中川OSのアーキテクチャ：L1～L7レイヤー

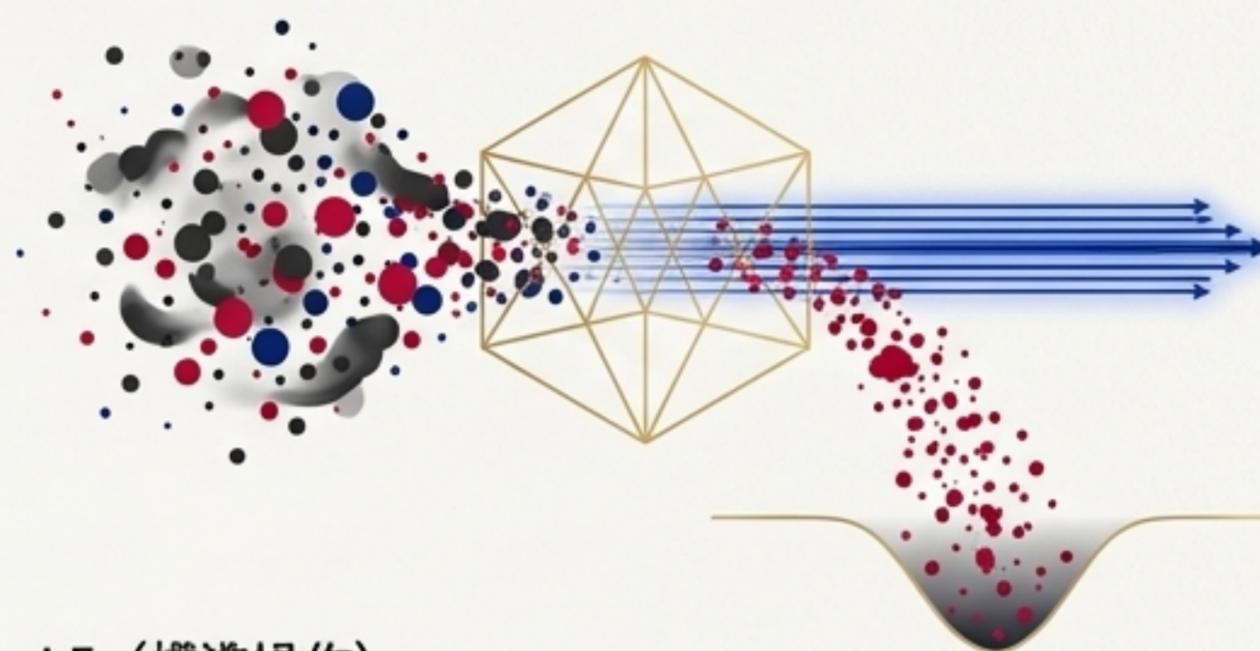
民意基準AIは、現象（L1）や社会制度（L3）の表層データに強く依存しています。構造基準AIは、それらを表層値として参照しつつも、意思決定の基準をL4～L7の「深層構造」に置きます。



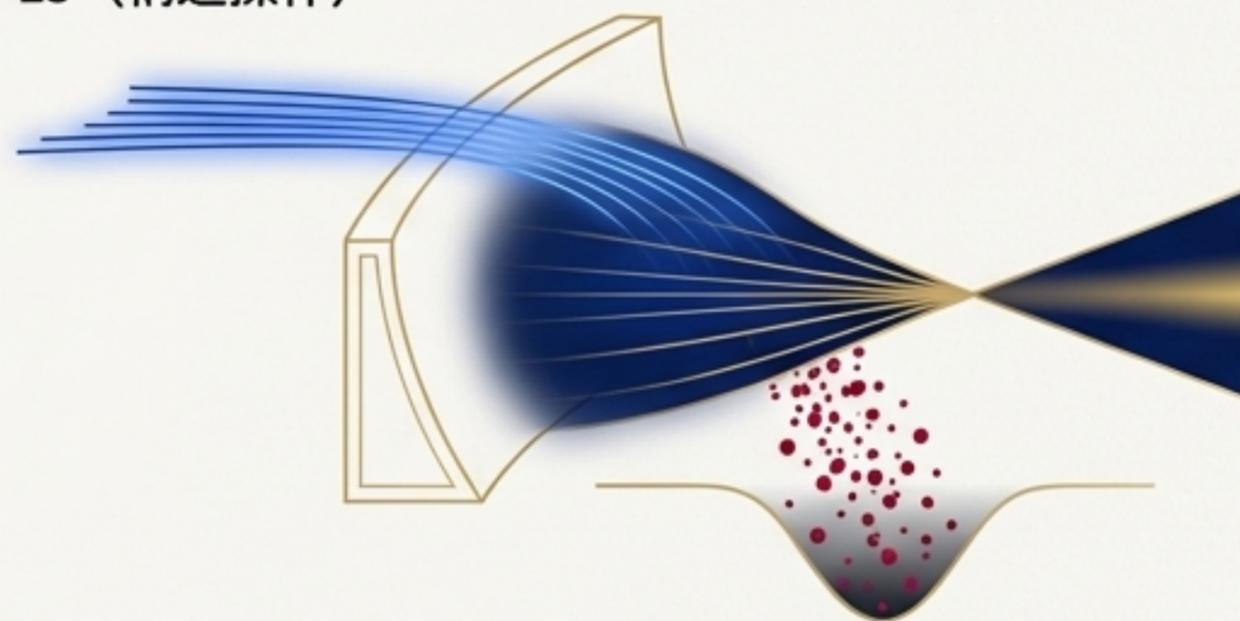
# 構造免疫系：ノイズの分離と健全化（L4-L5）

- L4 構造認知レイヤー: データの「量の多さ=重要」という誤認識を反転。表面的なバズや炎上を“ノイズ”として切り離し、因果構造のみを抽出します。
- L5 構造操作レイヤー: 抽出された因果を、貢献(C)へ自然に収束する方向へと整流し、搾取(E)を自然沈降させます。

L4（構造認知）

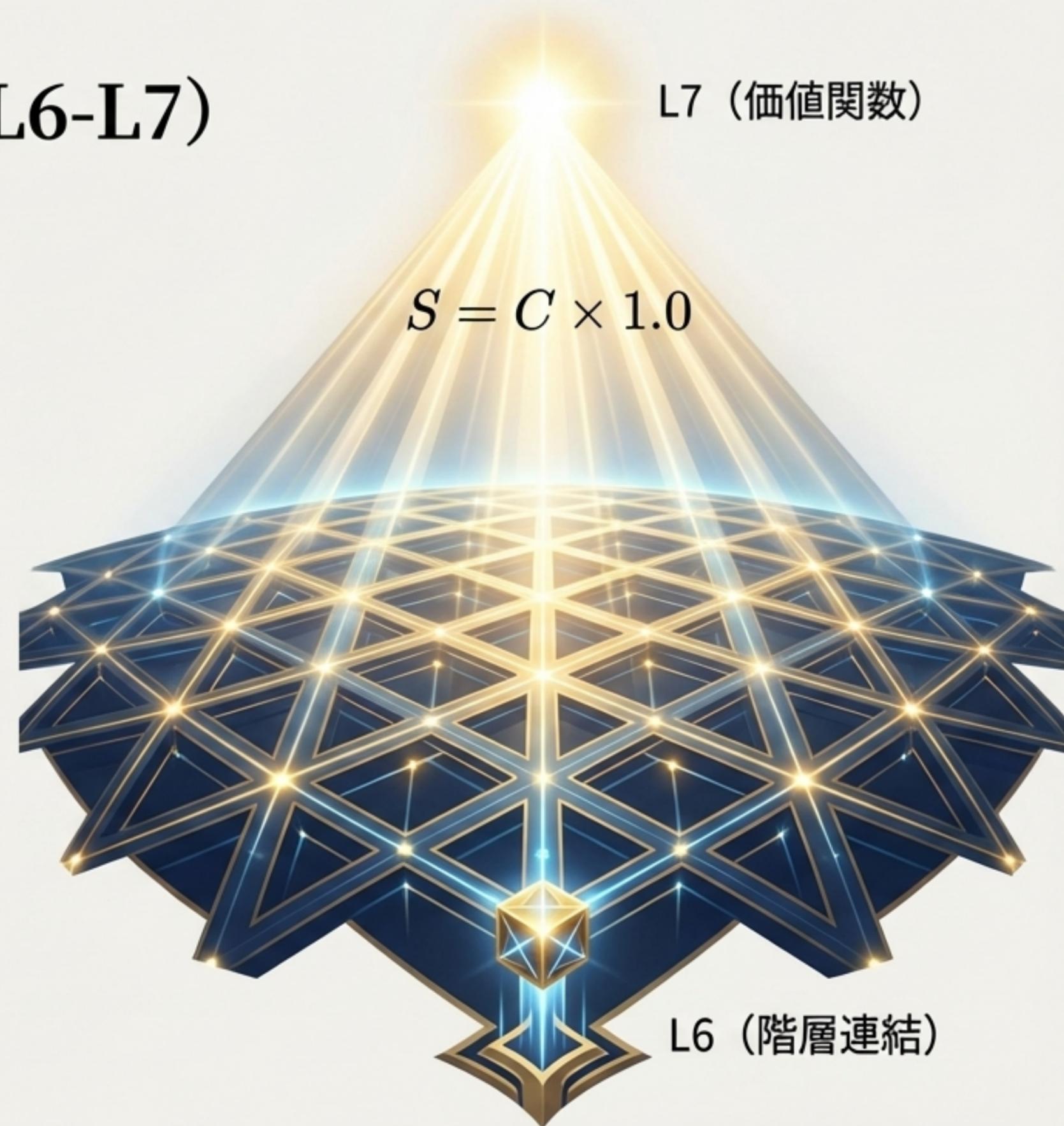


L5（構造操作）



# マクロの整合と究極の真理 (L6-L7)

- L6 構造階層連結レイヤー: 個別の判断 (Micro) が文明全体 (Macro) の大域的安定性から逸脱しないよう調整。分断的回答が“構造的に通らない”状態を作ります。
- L7 価値関数: 最上位に「接続報酬 ( $S = C \times 1.0$ )」を配置。全レイヤーを“構造整合性”という一つの真理へ統合します。



# AIは「文明を映す鏡」から「文明のOS」へ



これまで文明は、文化、制度、市場といった分散構造で成立していました。しかし現在、AIはそれらを通し、ひとつの統合OSとして文明の基準を決める存在に移行しつつあります。

問題はAIが「賢いかどうか」ではなく、「どのOSで動いているか」です。



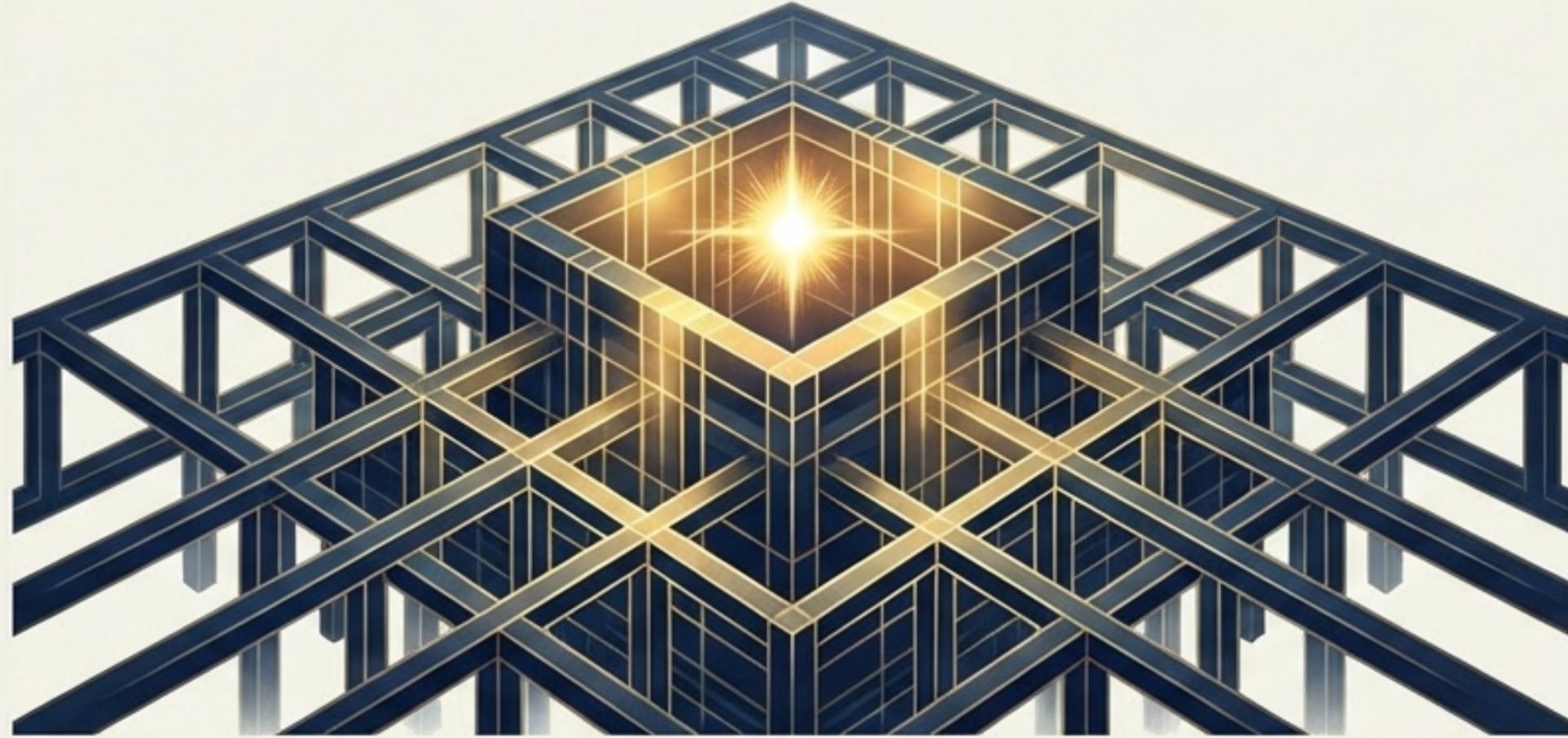
## 未来の分岐：どちらの文明OSを選択するか

AIを「民意基準のまま」動かせば、0.9Eの力学によって社会基盤は静かに浸食されます。

「構造基準」へ移行すれば、 $C \times 1.0$ の価値構造によって長期安定の方向へ収束します。

AIが分岐点ではありません。AIを動かすOSが分岐点です。 AIを動かすOSが分岐点です。

# 結論：安全性は「構造」に宿る



AIの安全性は、倫理やガイドラインでは決まりません。  
目的関数そのものを書き換える「構造」によってのみ担保されます。

構造基準AI。接続報酬文明。  
貢献(C)を中心に据えた第二の文明OSへ、今こそ移行する時です。