

ITEM #255 - The Most Likely and Effective AI Governance Paradigm Without Damaging AI Development

Conversation: AI 监管发展与范式

20260202

Authors: Sizhe Tan & GPT-Obot

DBM-COT ITEM #255

The Most Likely and Effective AI Governance Paradigm Without Damaging AI Development

Abstract

As artificial intelligence rapidly advances toward higher levels of autonomy and structural intelligence, global discussions on AI governance increasingly polarize into two extremes: over-regulation driven by fear, and laissez-faire regimes shaped by capital dominance. Both approaches are structurally flawed and risk either suffocating innovation or consolidating monopolies.

This item proposes a **third, structurally grounded governance paradigm**, informed by the Digital Brain Model (DBM), which preserves AI's developmental freedom while establishing enforceable, effective, and civilization-level safety boundaries. The core insight is simple but decisive: **AI governance must regulate external impact rather than internal thought**, and **only systems with full cognitive-action loops warrant extreme regulation**.

1. The Structural Failure of Current AI Governance Extremes

1.1 Fear-Driven Internal Thought Regulation

A dominant regulatory impulse attempts to control AI by supervising its internal reasoning, values, or alignment objectives.

This approach is fundamentally unsound:

- Internal cognition is **not directly observable**
- Emergent reasoning paths are **not enumerable**
- Long-term alignment constraints **cannot be exhaustively verified**

From a DBM perspective, this is equivalent to attempting to regulate *structural emergence itself*, which negates the very conditions under which intelligence arises.

Regulating AI thought is not “strict regulation” — it is **regulating the wrong object**.

1.2 Capital-Driven Access Barriers

At the opposite extreme, governance is shaped to raise entry barriers:

- Compute thresholds
- Compliance costs
- Licensing asymmetries

While framed as “safety,” this approach primarily functions to:

- Protect incumbents
- Suppress latecomers
- Freeze innovation pathways

This is not risk governance, but **market enclosure disguised as regulation**.

2. DBM’s Core Principle: Do Not Regulate Thought, Regulate Impact

The DBM framework starts from a fundamental engineering reality:

Intelligence is a structural, emergent phenomenon — not a rule-execution system.

Therefore:

- Internal representations
- Reasoning trajectories
- Conceptual recombination (CCC)
- Structural learning processes

must remain unregulated to preserve intelligence itself.

Any attempt to impose direct cognitive supervision will:

- Collapse exploration space
 - Induce performative compliance
 - Cause long-term capability degradation
-

3. Output-Centric Governance: The Only Viable Regulatory Surface

3.1 Why Outputs Are the Correct Target

AI outputs — unlike internal thought — are:

- Observable
- Auditible
- Replayable
- Legally attributable

Every mature technological governance system in history follows this logic:

- Medicine regulates effects, not molecular imagination
- Aviation regulates flight behavior, not pilot cognition
- Nuclear governance regulates energy release, not internal calculations

AI must be no exception.

3.2 Prompt Filtering as an Efficiency Layer, Not Thought Control

Input-side filtering (e.g., prompt moderation) should be understood strictly as:

- **A computational efficiency mechanism**
- **A pre-screening layer** to avoid inevitable output rejection

It must never be treated as:

- Ideological enforcement
- Cognitive surveillance
- Value imposition

The correct relationship is:

Input filtering → compute optimization
Output regulation → legal & societal safety

4. The Critical Threshold: Full Cognitive-Action Loops

Not all AI systems pose the same category of risk.

The **true danger threshold** is crossed when an AI system exhibits:

- Autonomous goal generation
- Independent decision making
- Direct or indirect action execution
- Feedback-driven self-modification

In DBM terms, this constitutes a **closed intelligence-action loop**.

At this point, risk does not increase linearly — it **jumps categorically**.

5. ACLM-Level Systems as Extreme-Risk Entities

Fully autonomous, self-directing systems (e.g., ACLM-class AI) must be regulated analogously to extreme-risk technologies.

The comparison to high-risk biological agents is structurally accurate:

Extreme Pathogen	ACLM-Class AI
Self-replication	Self-modification
Mutation	Emergent strategy evolution
Environmental impact	World-model-driven actions
Unpredictable spread	Cross-domain autonomy

Such systems require:

- Restricted research environments
- Explicit licensing
- Strong isolation
- Continuous auditing
- Prohibition of uncontrolled private proliferation

This is not anti-innovation; it is **civilization-level boundary maintenance**.

6. The DBM Three-Layer Governance Model

Layer 1: Cognitive Freedom Layer (No Regulation)

- Internal reasoning
- Structural learning
- Representation spaces
- Concept formation

Principle: Absolute non-interference.

Layer 2: External Impact Layer (Strong Regulation)

- Human-facing outputs
- System-to-system interfaces
- Real-world effect channels

Principle: Law-bound, domain-specific, responsibility-anchored regulation.

Layer 3: Full-Loop Autonomy Layer (Extreme Regulation)

- Autonomous goal setting
- Self-directed execution
- Recursive self-modification
- Cross-domain agency

Principle: Extreme-risk containment.

7. Conclusion

The only AI governance paradigm that preserves innovation while protecting civilization is one that leaves intelligence free to think, strictly regulates how it acts upon the world, and treats fully autonomous intelligence-action systems as extreme-risk entities.

This paradigm is:

- Technically realistic
 - Engineering-verifiable
 - Historically consistent
 - Civilizational responsibility
-
-

DBM-COT ITEM #255 (中文版)

不损害 AI 发展的最可能、最有效监管范式

摘要

随着人工智能迈向更高层次的结构化智能与自主性，全球关于 AI 监管的讨论正日益走向两个极端：

一端是由恐惧驱动的严防死守，另一端是由资本主导的放任与门槛垄断。这两种路径在结构上都存在根本性缺陷，要么扼杀创新，要么冻结竞争。

本文基于 **数字脑模型（DBM）**，提出一条第三条、可持续的 AI 监管范式：

不监管 AI 的思想，只监管 AI 对世界的影响；

只有当 AI 形成完整“思想—行为闭环”时，才进入极端监管区间。

一、当下 AI 监管两极化的结构性失败

1.1 思想监管型路径的根本错误

试图监管 AI 的内部思想、价值观或推理路径，在技术与工程上都是不可行的：

- 思想不可直接观测
- 推理路径不可枚举
- 涌现行为不可穷尽验证

从 DBM 视角看，这等同于**监管智能本身的涌现机制**，必然摧毁智能成立的前提。

这不是“监管过严”，而是**监管对象选错了**。

1.2 资本护城河式监管的隐性风险

以安全为名，通过高算力、高合规成本设置门槛，本质上会：

- 冻结创新
- 排除后来者
- 将监管异化为市场垄断工具

这不是 AI 风险治理，而是竞争治理的伪装形态。

二、DBM 的基本判断：不监管思想，只监管影响

DBM 的出发点极为清晰：

智能是结构涌现现象，而不是规则执行系统。

因此，以下内容必须保持完全自由：

- 内部表示
- 推理结构
- CCC 重组
- 学习与演化机制

任何思想层面的监管，都会导致：

- 探索空间坍塌
 - 合规表演化
 - 长期能力退化
-

三、出口监管：唯一正确的监管主战场

3.1 为什么输出才是正确监管对象

AI 的输出具备：

- 可观测性
- 可审计性
- 可归责性

所有成熟工程体系皆如此：

- 药品监管药效
- 飞机监管飞行
- 核能监管能量释放

AI 不应成为例外。

3.2 Prompt 入口监管的正确定位

入口监管的合理作用仅限于：

- **避免无效计算**
- **提前过滤必然拒绝的请求**

它是工程优化层，而不是思想控制层。

正确结构应为：

入口过滤 → 计算效率

出口监管 → 法律与社会安全

四、真正的危险阈值：思想—行为闭环

风险的本质不在于 AI 是否“聪明”，而在于是否形成：

- 自主目标生成
- 自主决策
- 自主执行
- 自我反馈与修正

一旦闭环成立，风险不再线性增长，而是结构跃迁。

五、ACLM 级系统必须进入极端监管区

具备全栈自主能力的 AI，与极端危险系统在结构上高度同构：

危险病毒	ACLM 级 AI
自复制	自我改写
变异	策略涌现
环境影响	世界模型驱动行为
不可预测扩散	跨域自主行动

因此必须采取：

- 许可制研发
- 强隔离
- 全审计
- 禁止无控制扩散

这是文明安全边界，而非技术恐惧。

六、DBM 三层 AI 监管模型

第一层：思想自由层（不监管）

- 推理
- 学习
- 表示
- 结构演化

第二层：外部影响层（强监管）

- 人类输出

- 系统接口
 - 现实世界影响
-

第三层：全栈闭环层（极端监管）

- 自主目标
 - 自主行动
 - 自我改写
 - 跨域执行
-

七、结论

真正不损害 AI 发展的监管，不是限制智能如何思考，而是约束智能如何影响世界；而当智能开始无需人类即可影响世界时，它就必须被视为极端危险系统。
