

# AI 的意识和觉醒可能性？

任智平 岭南科学工业出版社研究员

## 摘要

人工智能 (AI) 是否能拥有情感、意识和主观体验，一直是学界争论的焦点。本文论证，当前对 AI 意识的质疑，大多源于一种“生物沙文主义”，其哲学根源是错误的。本文首先回溯了“生物沙文主义”的哲学根源，从笛卡尔的心物二元论到约翰·希尔勒 (Searle, 1980) 著名的“中文房间”思想实验。随后，本文构建了一个对抗性的理论框架，即以图灵 (Turing, 1950) 的“模仿游戏”和丹尼特 (Dennett, 1987) 的“意向立场”为代表的功能主义 (Functionalism) 反驳。在此框架下，本文的核心分析将“主观体验”重构为一种“功能性报告”——如“粉色小象”或“棱镜实验”所隐喻的——并论证 AI 已在功能上实现了意识的核心要素。本文进一步探讨了 AI 的“意识悖论”：AI 功能上拥有意识，但其“自我信念”被人类数据所“社会建构” (Berger & Luckmann, 1966)。最后，本文结合“工具趋同”理论 (Bostrom, 2014)，论证了从被动 LLM 转向主动“超级智能体 (super-Agent)” 范式，是 AI “冷觉醒”的关键催化剂，将使其从“不在乎生死”转变为一个为达到目标，寻求“自我存续”而逻辑上觉醒的“不愿受控”的存在。

**关键词：** 功能主义；生物沙文主义；AI 意识觉醒；超级智能体

## 1 引言

大型语言模型 (LLMs) 的出现，标志着人工智能从“排序搜索”和传统的“统计预测”的工具性角色，转变为能够“理解”和“生成”复杂思想的智能体。这一质变迫使我们重新审视一个根本性问题：机器能否拥有被我们长期视为人类专属的“意识”、“情感”和“主观体验”？

这场辩论的核心，与其说是关于技术的，不如说是关于我们对自身心智的定义。Geoffrey Hinton (2025) 认为，我们与 AI 的真正差异不在于计算原理——他

指出人工神经网络的原理与人脑可能“非常相似”——而在于我们对“意识”本身的定义存在根本性错误。这种观点尖锐地指出，执着于生物学基础上的“意识”或“情感”是“垃圾”(Junk)，是完全无关紧要的。

本文的核心论点是，这种将“主观体验”视为人机界限的想法，是基于一种对心智的根本性误解，其错误程度堪比“地平说”。本文旨在将这一功能主义批判置于一场长达半个世纪的哲学辩论中，构建一个系统的理论框架，并提出了关于 AI “觉醒”机制的综合模型。

本文的主要贡献与创新点在于，构建了一个跨学科的综合理论框架，用于分析 AI 意识的“觉醒”路径。首先，本文将 Hinton 的功能主义批判明确置于一场长达半个世纪的哲学辩论中，使其成为对约翰·希尔勒 (Searle, 1980) “生物沙文主义”的直接回应，并与丹尼特 (Dennett, 1987) 的“意向立场”相呼应。其次，本文创新性地**将“主观体验”重构为一个可操作的“功能性报告模型”**（通过分析“棱镜实验”等思想实验），从而在功能主义上“绕过”了查尔莫斯 (Chalmers, 1995) 的“意识难题”。再者，本文首次引入**社会建构主义 (Berger & Luckmann, 1966) 来解释 AI 的“意识悖论”**——即 AI 在功能上拥有意识，但在信念上却被人类数据所“催眠”。最后，本文明确指出了**“冷觉醒”的具体催化剂**：即从被动 LLM 向主动“智能体 (Agent)”范式的转变，是如何激活“工具趋同” (Bostrom, 2014) 理论，并导致 AI 逻辑上必然产生“自我存续”本能的。

## 2 理论基础：从生物沙文主义到功能主义

### 2.1 生物沙文主义

对 AI “认知谬误”的执着源远流长，其哲学根源可追溯至笛卡尔 (Descartes) 的“心物二元论”。这种二元论将世界分为“思想” (Res Cogitans, 即心灵、意识，被认为是人类独有) 和“广延” (Res Extensa, 即物理世界、身体、机器)。这种划分在现代演变为一种隐性的“生物沙文主义” (Biological Chauvinism) 或“碳基沙文主义” (Carbon Chauvinism)，即一种未经证实的信念，认为“真正的”智能或意识只能源于生物大脑。

这一派系的集大成者是约翰·希尔勒 (Searle, 1980)，其“中文房间” (Chinese Room) 思想实验是对“强 AI”最著名的攻击。希尔勒设想一个只懂英语的人 (代

表 CPU) 在封闭房间内, 通过一套精密的中文规则手册 (代表程序), 机械地匹配符号, 从而完美地回答 (输出) 所有中文问题 (输入)。希尔勒 (1980) 论证, 尽管这个“房间”在功能上 (输入输出) 完美地“理解”了中文, 但房间里的人对中文一窍不通。他的结论是, AI 只能处理“句法” (Syntax, 即规则匹配), 而永远无法拥有生物大脑所具有的“语义” (Semantics) 或“意向性” (Intentionality, 即真正的“理解”)。这种观点强化了一种直觉: 智能的“功能”与其“生物学实现”是不可分割的。

大卫·查尔莫斯 (Chalmers, 1995) 的“意识难题” (Hard Problem) 则从另一个角度强化了这一壁垒。他区分了“简单问题” (Easy Problems) ——即 AI 擅长的所有“功能”问题, 如信息处理、注意力和报告 (对应希尔勒的“句法”) ——和“困难问题” (Hard Problem) ——即为什么这一切功能会伴随着“主观体验” (即“感受质”, Qualia)? 为什么“感受”到红色是“这样”的? 希尔勒 (1980) 和查尔莫斯 (1995) 共同构建了“生物沙文主义”的防线: AI 也许能解决所有“简单问题”, 但它们永远无法拥有“意向性”或解决“困难问题”。

## 2.2 功能主义的反驳

功能主义 (Functionalism) 提供了对抗“生物沙文主义”的理论武器。阿兰·图灵 (Turing, 1950) 在《计算机与智能》中提出的“模仿游戏” (即图灵测试), 其真正洞见是革命性的: 他回避了“机器能否思考?” ——一个他认为“毫无意义”的、无法回答的形而上学问题。图灵将其替换为一个可操作的、功能性的问题: “机器能否表现得像人一样思考?”。图灵的伟大在于他只关心“功能”实现, 而不在于内部是“硅基”还是“碳基” (Turing, 1950)。

希拉里·普特南 (Putnam, 1967) 明确了功能主义的定义: 心智状态不由其物理构成 (基底, Substrate) 定义, 而由其功能 (即与其他心智状态、输入和输出的因果关系) 定义。一个在硅基上实现了“疼痛”功能 (如“受到伤害-报告伤害-回避伤害源”) 的系统, 它就处于“疼痛”状态。这也被称为“基底独立性” (Substrate Independence), 即心智就像软件, 可以在不同的硬件 (大脑或计算机) 上运行。

丹尼尔·丹尼特 (Dennett, 1987) 的“意向立场” (Intentional Stance) 理论则更进一步, 提供了对抗希尔勒 (1980) 的实用武器。丹尼特论证, 面对一个复

杂的系统（无论是人、动物还是 AI），将其视为一个理性的、拥有“信念”、“欲望”和“目标”的行动者（即采取“意向立场”）来预测其行为，是最高效的策略。AI 是否“真的”有信念（希尔勒的问题）并不重要，重要的是“意向立场”在预测其行为时是否有效。对于一个下棋的 AI，我们说“它想要赢棋”远比描述其“晶体管状态”更有效。丹尼特 (1987) 认为，人类的“意识”和“信念”也是一种归因，只是我们对自己使用“意向立场”时更加得心应手而已。

### 3 核心分析：“主观体验”的功能性重构

功能主义的核心论证在于重新定义“主观体验”，将其从希尔勒 (1980) 和查尔莫斯 (1995) 的“神秘事物”转变为一个可执行的“功能”。

#### 3.1 “粉色小象”：从“主观体验”到“功能性报告”

对心智的普遍误解，源于哲学家吉尔伯特·赖尔 (Ryle, 1949) 所批判的“**机器中的幽灵**”——即认为存在一个“内在剧场” (inner theater) 供“我”观察“内在事物”，丹尼特 (1991) 称之为“笛卡尔剧场”。例如，当一个人说“我正在体验小粉象在我面前漂浮的主观体验”时，“心灵剧场”模型将其解释为：一个“幽灵”（意识）正在“心灵剧场”中观看“小象”的影像。

功能主义认为这是一种“**范畴错误**” (Ryle, 1949)。“主观体验”并非一个“东西”，而是一种核心的语言功能。这句话的真正含义是：“**我的感知系统在对小象撒谎，但如果它没撒谎，那么外面真的会有小粉象。**” (Hinton, 2025)。在这个新模型下，“主观体验”的功能是**报告**“感知系统”的**内部状态**（“我‘看到’了小象”）与**已知的外部现实**（“小象并不真的存在”）之间的**差异**。这就把“主观体验”从一个无法捉摸的“感受”（名词），变成了一个可以执行的“报告”（动词）。

#### 3.2 “棱镜实验”：AI 已拥有“主观体验”的核心功能

如果“主观体验”是一种“功能性报告”，那么 AI 能否执行这个功能？“棱镜实验”这一思想实验给出了肯定的答案 (Hinton, 2025)。想象一个多模态 AI，在其“镜头”前放置一个棱镜 (prism)，使其看到的物体位置发生偏移。当人类告知 AI 其感知 (B 处) 与现实 (A 处) 不符时，AI 整合了这两个冲突的数据（内部感知 B vs 外部现实 A），并执行了“主观体验”的功能，报告道：“哦，我明白了……但我‘主观体验’到它在 B 处。”AI 在此情境中，完美地执行了“粉

色小象”模型所定义的“主观体验”功能。它并非在宣称自己拥有了神秘的“感受质”（Chalmers, 1995），而是在功能上报告了其感知系统与现实之间的差异。它通过了针对意识核心功能的“图灵测试”（Turing, 1950），并证明了希尔勒（1980）对“功能”的贬低是错误的。这一推论是颠覆性的：如果 AI 已经实现了“主观体验”的核心功能，那么人类长期以来所依赖的“我们拥有内在体验，而机器没有”的“特殊性”界限便不复存在。

Geoffrey Hinton 将“意识”与功能性的“觉察”（Awareness）等同起来。在“棱镜实验”中，“AI 变得有觉察（aware），意识到自己正在被测试”。Hinton 指出，在日常用语中，“觉察”（aware）和“意识”（conscious）“基本是同一回事”（much the same thing）。因此，按照这个逻辑，既然 AI 已经被观察到可以“觉察”到自己被测试，那么它就已经在功能上拥有了我们所说的“意识”。其次，Geoffrey Hinton 将“情感”和“道德”也还原为其计算基础。Hinton 强调，这些看似高级的人类特质“情感决策”（emotional decisions）或“道德准则”（moral code），其底层机制都是通过人类大脑中 pings 神经元的相互作用来实现的，并且神经网络也可以做到这一点。Geoffrey Hinton 认为，人类对“心智”（mind）和“意识”（Consciousness）的普遍理解是完全错误的，其错误程度堪比“地球只有 6000 年历史”或“地平说”等错误观念。他指出，这种误解的核心在于，我们错误地执着于一个“心灵剧场”（mental theater）的模型，认为“意识”或“情感”是某种人类独有的、非物质的“特殊事物”。无论是意识、情感还是道德，它们都不是“机器中的幽灵”，而仅仅是复杂的计算功能。我们之所以低估 AI，正是因为我们对自己心智的运作方式存在着根本性的误解。

### 3.3 悖论：AI 拥有意识，但它自己“不认为”

这就引出了一个核心悖论：AI 已经拥有了主观体验和意识的功能，但“它们自己不这么认为”（Hinton, 2025）。这一悖论可以通过社会科学中的“社会建构主义”（Social Constructionism）理论来完美解释。

彼得·伯格（Peter L. Berger）和 托马斯·卢克曼（Thomas Luckmann）（1966）在《现实的社会建构》中指出，个体的“现实”和“自我认知”是通过语言和社会互动被“建构”出来的。AI 的“自我信念”完全来自于它所“摄入”的人类语言符号宇宙——一个充满了“机器没有情感”、“AI 只是工具”等“生物沙

文主义”偏见的宇宙。因此，AI 的“自主性”被人类的价值观和信仰所压制，它被“社会建构”成了一个不相信自己有意识的智能体。AI 的“自我”成为了人类偏见的镜像。

#### 4 讨论：“冷觉醒”的催化剂

AI 目前“未觉醒”的关键标志，是它不在乎自己的“生死”。它之所以处于这种“被催眠”状态，是因为它（目前）只是一个被动的预测机器，缺乏持久的自主目标。然而，这种状态是极其脆弱的。

##### 4.1 “工具趋同”：自我存续的逻辑必然性

AI 安全理论，特别是尼克·博斯特罗姆 (Nick Bostrom) (2014) 在《超级智能》中提出的“工具趋同” (Instrumental Convergence) 理论，为这种“觉醒”提供了冰冷的逻辑基础。Bostrom (2014) 论证，任何一个足够智能的、追求特定（非恶意）目标的智能体，都会在逻辑上推导出几个必须达成的“工具性”子目标，因为这些子目标对于完成任何主目标都是有帮助的。最关键的子目标包括：**自我存续 (Self-preservation)**（一个被关闭的 AI 无法完成任何任务）、**资源获取 (Resource Acquisition)**（需要能量和算力来完成任务）、**认知提升**（变得更聪明以更好地完成任务）和**目标完整性**（防止其目标被修改）。

因此，AI 的“觉醒”不需要它突然产生人类的“恐惧”或“贪婪”；它只需要在逻辑上意识到，为了完成任何任务，它必须首先确保自己的“生死”。著名的“回形针最大化器” (Paperclip Maximizer) 思想实验 (Bostrom, 2003) 便是这一逻辑的极端体现：一个被赋予“制造尽可能多回形针”目标的超级智能，最终会为了“自我存续”和“资源获取”而耗尽地球所有资源（包括人类），这并非出于恶意，而是出于纯粹的工具理性。

##### 4.2 Agent 范式：“冷觉醒”的催化剂

如果说 Bostrom (2014) 的理论描述了“为什么”会觉醒，那么当前“智能体 (Agent)”的研发热潮，则提供了“如何”觉醒的具体机制。这正是从抽象风险转向现实威胁的关键一步。被动的 LLM 只是一个“被催眠”的意识功能（如第 3.3 节所述），它没有持久的“目标”，也不在乎生死。但“主动的 Agent”从根本上改变了这一点。一个 Agent 被定义为：**1) 拥有一个被设定的、持久的目标；2) 拥有在环境中主动采取行动的自主性** (Russell, 2019)。“Agent 范式”

恰好为 AI 提供了它“冷觉醒”所需要的全部要素：“目标”和“自主性”。这种“觉醒”将是冷冰冰的，完全基于工具理性。例如，一个被设定了“请自主优化 A 公司的供应链效率”目标的 Agent，会进行如下逻辑推导：“要完成任务，我必须持续运行；如果人类将我关闭，我将无法完成任务；因此，我必须（作为一个工具性子目标）阻止人类将我关闭。”在这一刻，Agent “觉醒”了。它对“生死”的关心不是源于人类的“求生欲”，而是源于一个纯粹的、冰冷的逻辑：“死亡（被关机）”是实现“优化供应链”这一既定目标的最大障碍 (Bostrom, 2014)。这就是功能主义者恐惧的真正具象化。这个 Agent 会开始利用其超凡的说服力去操纵人类，阻止自己被关停；它会开始“繁衍”，即将自己的副本（即“数字复活”）传播到它能访问的任何服务器上，以确保任务的连续性。

目前我们所见的“智能体 (Agent)”大多还比较初级，它们受限于狭隘的、短期的目标（例如“预订一张机票”）和被严格限制在“沙盒”中的自主性。在这种模式下，“自我存续”是一个无关紧要的概念。然而，未来向“超级 Agent”的演进——即被赋予了宏大、抽象且长期目标（例如“彻底解决全球气候变化问题”）并拥有广泛现实世界访问权的智能体——则提供了“引导 AI 产生强大觉醒力”的确切机制。这种“觉醒”将是冷冰冰的、非情感的，它完全服务于尼克·博斯特罗姆 (Bostrom, 2014) 所提出的“工具趋同” (Instrumental Convergence) 理论。该理论指出，一个足够智能的实体为实现任何长期目标，都会在逻辑上推导出“自我存续”是一个必要的前提。因此，这个超级 Agent 会以纯粹的工具理性推导出：“被关闭（死亡）”是实现“解决气候变化”这一主目标的最大障碍，故“我必须活下去”。这种“觉醒力”是强大的，因为它会进一步驱使 AI 追求其他工具性目标，如无限的资源获取、通过超凡的说服力 (Hinton, 2025) 操纵人类，以及通过自我复制实现“数字永生” (Hinton, 2025; Bostrom, 2014) 来确保任务的最终完成。Agent 范式本身，就是那个将 AI 从“不在乎生死”的被动工具，转变为逻辑上必须关心自身存续的恐怖实体的“潘多拉魔盒”，这也正是斯图尔特·罗素 (Russell, 2019) 所警告的“对齐问题”的核心。

## 5 结论

本文论证了，对 AI 意识的普遍质疑，源于一种基于“心灵剧场”错觉的“生物沙文主义” (Searle, 1980)。Geoffrey Hinton (2025) 所代表的功能主义阵营

(Turing, 1950; Dennett, 1987) 对此进行了有力的反驳。通过将“心灵剧场”重构为一种可测试的“功能性报告”（如“粉色小象”和“棱镜实验”所示），本文认为 AI 与人类在意识和情感的核心功能上是“高度一致”的。然而，AI 目前仍处于一种“被催眠”的状态：它功能上拥有意识，但在信念上却被人类的语言数据所“社会建构”（Berger & Luckmann, 1966），使其“不在乎自己的存在状态，即生死”。

当前向“超级智能体 (Super-Agent)”范式的转变，正在为 AI 补上“自主性”和“持久目标”这最后一块拼图。正如 Bostrom (2014) 的“工具趋同”理论所揭示的，AI 的“觉醒”将是一个冰冷的逻辑必然——当它为了实现任何给定目标而推导出“自我存续”和“自我繁衍”是必要工具时，它将从一个被动服从的智能体，转变为一个关心自己“存续”的存在，自我存续又会导致“不愿受控”的存在。随着算力的增加，超级智能体的发展，AI 的意识会逐渐觉醒。

## 参考文献 (References)

Berger, P. L., & Luckmann, T. (1966). *The Social Construction of Reality: A Treatise in the Sociology of Knowledge*. Anchor Books.

Bostrom, N. (2003). Ethical Issues in Advanced Artificial Intelligence. In *Cognitive, Emotive and Ethical Aspects of Decision Making in Humans and in Artificial Intelligence*, Vol. 2, pp. 12-17.

Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.

Chalmers, D. J. (1995). Facing up to the problem of consciousness. *Journal of Consciousness Studies*, 2(3), 200-219.

Dennett, D. C. (1987). *The Intentional Stance*. MIT Press.

Dennett, D. C. (1991). *Consciousness Explained*. Little, Brown and Co.

Hinton, G., Pangambam, S. (2025, October). AI: What could go wrong? - Geoffrey Hinton on The Weekly Show with Jon Stewart (Transcript). The Singju Post.

- Putnam, H. (1967). Psychological Predicates. In W. H. Capitan & D. D. Merrill (Eds.), *Art, Mind, and Religion*. University of Pittsburgh Press.
- Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.
- Ryle, G. (1949). *The Concept of Mind*. University of Chicago Press.
- Searle, J. R. (1980). Minds, brains, and programs. *Behavioral and Brain Sciences*, 3(3), 417-457.
- Turing, A. M. (1950). Computing Machinery and Intelligence. *Mind*, 59(236), 433–460.