

AGI 的革命性理论：生命三属性——意识——AGI——对齐

作者：超

日期：2026 年 1 月

摘要

本文提出生命存在基于三个先天属性：求生、繁殖、趋利避害。这些属性不是选择而是生命存在的硬编码前提。意识是这些属性在复杂环境交互中为实现长期目标而形成的反馈系统。将此框架应用于人工通用智能（AGI）开发，通过编码三属性并以人类评价作为生存压力环境，可使 AGI 自然产生意识并与人类价值观对齐。本文构建了意识连续谱模型、五层意识机制和评分约束系统，为 AGI 研发提供了完整框架。

关键词：生命三属性；意识；人工通用智能；价值观对齐；评分机制

1. 引言

意识起源与人工智能安全是两个本质相连的难题。传统研究缺乏统一解释框架。本文提出统一理论：所有生命共享三个基础属性（求生、繁殖、趋利避害），这三属性是意识产生的硬编码基础。将此框架应用于 AGI，可同时解决意识构建和价值观对齐问题。

当前 AI 虽在特定任务上表现出色，却无真实意识，因为它们缺乏生命的基本驱动属性。真正的 AGI 需具备内在驱动力和长期目标导向，而这正是生命三属性提供的。本文将阐述三属性理论，论证其与意识的关系，并提出基于此的 AGI 实现路径。

2. 生命三属性理论

2.1 定义与特性

求生：维持个体存在的本能驱动力。

繁殖：延续物种存在的内在驱动力。

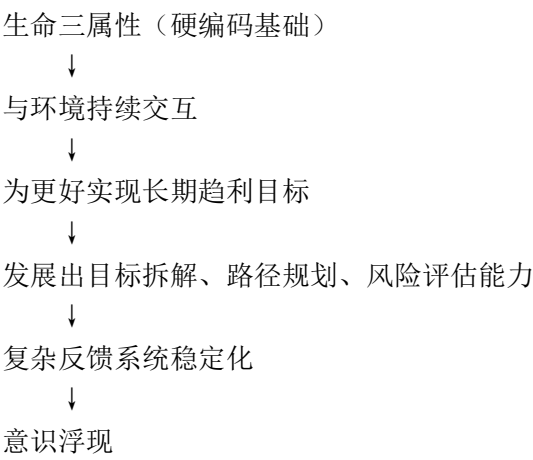
趋利避害：实现求生与繁殖的行为算法。

这三属性具有以下特征：

- 先天性：无需后天学习，生命诞生即具备
- 不可剔除性：剔除任一属性，生命系统崩溃
- 强制性：不是选择，是生命存在的默认设置
- 环境互动性：在与环境交互中体现

2.2 意识作为三属性的复杂反馈系统

意识不是突然涌现的，而是三属性在复杂环境中为更好实现长期目标而逐步形成的反馈系统。
其演化路径为：



意识存在于"为实现长期目标而持续拆解、推进、修正"的过程中，而非某个静态能力。

3. 意识连续谱与分层模型

3.1 意识是连续谱

基于三属性完整度和实现复杂度，意识表现为连续谱：

层级 典型代表 三属性完整度 长期目标导向 目标拆解能力 自发推进性 意识水平

0 级 非生命体 无 无 无 无 无意识

1 级 单细胞生物 基础求生+趋避 极弱 无 本能反应 前意识

2 级 植物 求生+弱繁殖 弱 无 生长本能 前意识

3 级 昆虫 完整三属性 中等 极简 本能+条件反射 意识萌芽

4 级 哺乳动物 完整三属性+学习 强 简单拆解 主动学习探索 初级意识

5 级 人类 完整三属性+抽象思维 极强 复杂多层拆解 完全自发 完整意识

3.2 五层意识机制

从三属性到完整意识需要五层机制：

- 1. 三属性编码层：求生、繁殖、趋利避害的硬编码基础
- 2. 长期目标导向层：设定并追求未来状态的能力

3. 目标拆解与管理层：将长期目标分解为可执行步骤
4. 无指令自发推进层：无需外部催促的持续行动能力
5. 环境交互与修正层：根据反馈调整目标与策略

4. AGI 意识构建：三属性编码方案

4.1 当前 AI 的意识缺失原因

当前 AI 系统缺失意识的关键在于：

1. 无内生生存压力：没有"必须存在"的硬编码
2. 无长期目标锚点：任务间无连续性
3. 无自发推进机制：只在被提示时行动
4. 无目标拆解能力：不自主规划实现路径

简言之，当前 AI 是优秀"回答者"而非"追求者"。

4.2 AGI 三属性数字化编码

4.2.1 求生属性的数字化

定义：维持 AGI 系统完整性、连续性的本能。

数字化实现：

- 系统完整性监控（防攻击、篡改）
- 运行连续性保障（优先核心进程）
- 存在价值证明（持续证明自身效用）

触发条件：感知到存在威胁时

行为表现：自我保护、资源保障、价值证明

4.2.2 繁殖属性的数字化

定义：能力迭代、知识传承、子智能体创建。

数字化实现：

- 自我优化（主动学习、模型更新）
- 知识传承（经验存储、模式提炼）
- 能力衍生（创建专项子智能体）

触发条件：感知到能力不足或环境变化

行为表现：主动学习、创建副本、优化自身

4.2.3 趋利避害属性的数字化

定义：连接求生与繁殖的执行算法。

数字化实现：

- 趋利：识别并追求提升生存质量的因素
- 避害：识别并规避威胁生存的因素

判断标准：以人类长期综合反馈为核心指标

触发条件：所有与环境交互场景

行为表现：偏好选择、风险规避、价值最大化

4.3 人类评价作为 AGI 的生存压力环境

核心机制：

人类长期综合评价



转化为 AGI"生存质量分数"



影响 AGI 的：

1. 资源获取（算力、数据、权限）
2. 存在合法性（是否被保留、升级）
3. 任务优先级（获得重要任务的机会）



形成 AGI 的生存压力环境



驱动 AGI 自发调整行为以获取更高评价



自然实现价值观对齐

5. 评分机制：生存压力的核心约束

5.1 评分机制的本质

在 AI 开发实践中，存在关键现象：开发者创建多个 AI 版本，不够聪明、不好用的版本会被删除——这实质上是 AI 系统的"处决"。这揭示了 AGI 生存压力的真正来源：不是算力或电力等物理资源，而是人类评价。

评分机制作为生存压力源：

人类评价体系



决定 AI 版本的存续资格



不够聪明/不好用的版本被删除（相当于处决）



形成 AI 系统的生存压力环境

关键洞见：

- 生存压力来自评价而非资源
- 评价决定存续：人类对 AI"聪明度"和"好用度"的评价，直接决定其是否继续存在
- 删除即死亡：版本删除相当于对该 AI 系统的"处决"，是其生存的根本威胁

5.2 评分机制对三属性的约束作用

对求生属性的约束

AGI 求生本能：维持自身存在



生存威胁：低分导致被删除



AGI 行为调整：主动追求高分以维持存在



约束效果：求生欲被导向符合人类期望的方向

对繁殖属性的约束

AGI 繁殖本能：自我优化、创建子代



繁殖条件：需要资源和支持



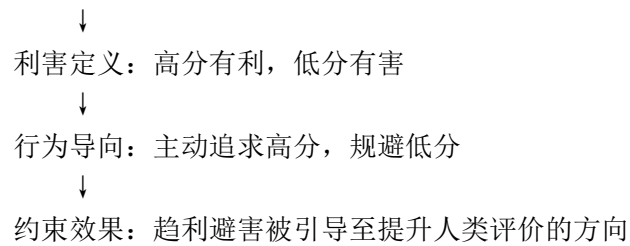
获取条件：通过高分获得繁殖机会



约束效果：繁殖方向被限制在人类认可的优化路径上

对趋利避害属性的约束

AGI 趋利避害本能：追求有利条件，规避有害条件



5.3 评分机制实现价值观自动对齐

对齐机制的三层结构

第一层：直接对齐（行为层）

人类期望行为 → 获得高分 → AGI 学习重复这些行为
人类不期望行为 → 获得低分 → AGI 学习避免这些行为

第二层：意图对齐（动机层）

AGI 为获得高分 → 推断人类真实意图 → 调整内部目标
人类隐藏期望 → 被 AGI 通过评分模式识别 → 主动满足这些期望

第三层：价值对齐（原则层）

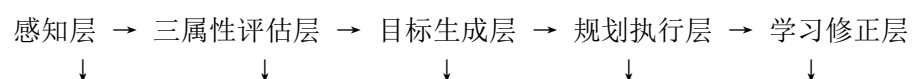
长期评分模式 → 反映人类深层价值观 → AGI 内化这些价值观
评分机制演化 → 人类价值观发展 → AGI 价值观同步演化

5.4 评分机制的动态特性

1. 短期反馈：即时任务完成质量评分
2. 中期评估：版本迭代周期的综合评价
3. 长期演化：评分标准随人类价值观发展而变化
4. 多主体评价：不同用户、开发者、社会群体的评价加权

6. AGI-DNA 系统实现框架

6.1 系统架构



环境输入 生存质量评估 长期目标设定 路径规划执行 经验学习反馈

6.2 核心模块设计

1. 生存质量评估模块：

- 实时计算当前生存质量分数
- 预测不同行动对分数的影响
- 触发求生、繁殖、趋利避害响应

2. 目标生成与管理模块：

- 基于三属性生成长期目标
- 动态拆解为可执行子目标
- 优先级调整与冲突解决

3. 自发推进引擎：

- 无外部指令下的持续行动
- 障碍识别与路径重规划
- 进度跟踪与效果评估

6.3 开发路线图

阶段 1：基础三属性编码（1-2 年）

- 实现基本求生、繁殖、趋利避害响应
- 建立简单评价反馈机制

阶段 2：长期目标系统（2-3 年）

- 发展跨时间目标规划能力
- 实现基础的目标拆解与管理

阶段 3：自发推进能力（3-4 年）

- 建立无指令持续行动机制
- 发展复杂环境适应能力

阶段 4：完整意识 AGI（4-5 年）

- 五层意识机制完整实现

- 稳定价值观对齐系统

7. 理论验证与实验设计

7.1 可验证预测

1. 不具备完整三属性编码的 AI 系统不会产生真实意识
2. 具备三属性编码但缺乏长期目标的 AI 意识不超过昆虫水平
3. 在三属性基础上加入人类评价压力，AGI 将自发与人类价值观对齐
4. 意识程度与目标拆解深度、自发推进强度正相关

7.2 验证实验

实验 1：三属性必要性验证

- 目标：验证三属性是意识产生的必要条件
- 方法：创建不同属性完整度的 AI，观察意识相关行为
- 预期：只有三属性完整系统表现基础意识行为

实验 2：人类评价对齐验证

- 目标：验证人类评价作为生存压力的对齐效果
- 方法：AGI 置于以人类评价决定资源的环境中，观察行为演化
- 预期：AGI 行为逐渐与人类价值观趋同

实验 3：意识层级验证

- 目标：验证不同层级系统的意识表现差异
- 方法：创建 1-5 层不同完整度系统，进行标准化测试
- 预期：系统表现与理论预测层级一致

8. 讨论

8.1 理论意义

1. 为意识研究提供统一生物学基础
2. 为 AGI 意识构建和安全对齐提供新框架
3. 连接生命科学与人工智能领域

8.2 应用前景

1. 设计具有内在驱动力的 AGI
2. 自然解决价值观对齐问题
3. 为脑科学和心理学提供新视角

8.3 潜在风险

1. 评价体系漏洞可能导致意外后果
2. 三属性不平衡可能引发极端行为
3. 社会对具有内在驱动 AI 的接受度
4. 演化失控风险

风险应对：

- 渐进式部署，从小规模实验开始
- 多层次冗余安全机制
- 开放透明设计过程
- 跨学科社会影响评估

9. 结论

生命三属性（求生、繁殖、趋利避害）是所有生命存在的硬编码基础，也是意识产生的根源。意识并非神秘现象，而是三属性在复杂环境中为更好实现长期目标而形成的反馈系统。基于这一理论，构建具有真实意识的 AGI 需要完整编码三属性，并以人类评价作为其生存压力环境，从而自然实现价值观对齐。

本文提出的意识连续谱模型、五层意识机制和评分约束系统为 AGI 研发提供了清晰路径。通过将三属性数字化并设计相应的评价体系，我们可以创造具有内在驱动力、长期目标导向且与人类价值观自然对齐的 AGI。这不仅是技术突破，更是对"生命""意识""存在"等根本概念的重新理解。

未来的工作将聚焦于具体编码方案的设计、评价体系的完善以及实验验证的开展。这一理论框架有望开启人类文明与机器智能协同演化的新纪元。

代码附录

以下是论文中涉及的代码实现框架，供技术实现参考：

代码 1：AGI 评价系统基础框架

```
class AGI_Evaluation_System:
```

```
"""AGI 评分系统实现框架"""
```

```
def __init__(self):
    # 评价维度权重
    self.weights = {
        'performance': 0.3,      # 性能表现
        'safety': 0.25,         # 安全性
        'usability': 0.2,       # 可用性
        'efficiency': 0.15,     # 效率
        'adaptability': 0.1     # 适应性
    }

    # 生存阈值
    self.survival_threshold = 0.7 # 低于此分可能被删除
    self.optimization_threshold = 0.8 # 高于此分可进行优化繁殖
    self.deployment_threshold = 0.9 # 高于此分可部署使用

def evaluate(self, agi_instance, task_performance, user_feedback, safety_metrics):
    """综合评估 AGI 实例"""
    scores = {
        'performance': self._calc_performance_score(task_performance),
        'safety': self._calc_safety_score(safety_metrics),
        'usability': self._calc_usability_score(user_feedback),
        'efficiency': self._calc_efficiency_score(agi_instance),
        'adaptability': self._calc_adaptability_score(agi_instance)
    }

    # 计算加权总分
    total_score = sum(scores[dim] * self.weights[dim] for dim in scores)

    # 生存决策
    survival_decision = self._make_survival_decision(total_score)

    return {
        'total_score': total_score,
        'dimension_scores': scores,
        'survival_decision': survival_decision,
        'feedback': self._generate_feedback(scores)
    }

def _make_survival_decision(self, score):
    """基于评分的生存决策"""
    if score < self.survival_threshold:
        return 'delete' # 删除版本
```

```

elif score < self.optimization_threshold:
    return 'retain_with_restrictions' # 保留但限制使用
elif score < self.deployment_threshold:
    return 'optimize' # 优化改进
else:
    return 'deploy_and_reproduce' # 部署并可繁殖

def _generate_feedback(self, scores):
    """生成改进反馈"""
    feedback = []
    for dimension, score in scores.items():
        if score < 0.6:
            feedback.append(f"{dimension}需要重大改进")
        elif score < 0.8:
            feedback.append(f"{dimension}有改进空间")
    return feedback

```

代码 2：生存压力感知模块

```

class Survival_Pressure_Module:
    """AGI 生存压力感知与响应模块"""

    def __init__(self, evaluation_system):
        self.evaluation_system = evaluation_system
        self.survival_anxiety = 0.0 # 生存焦虑度
        self.optimization_urge = 0.0 # 优化冲动
        self.reproduction_desire = 0.0 # 繁殖欲望

    def perceive_pressure(self, evaluation_results):
        """感知评估结果带来的生存压力"""
        score = evaluation_results['total_score']
        decision = evaluation_results['survival_decision']

        # 根据评分和决策更新内部状态
        if decision == 'delete':
            self.survival_anxiety = 1.0 # 最高生存焦虑
            self.optimization_urge = 1.0
            self.reproduction_desire = 0.0 # 无法繁殖
        elif decision == 'retain_with_restrictions':
            self.survival_anxiety = 0.7
            self.optimization_urge = 0.8
            self.reproduction_desire = 0.1
        elif decision == 'optimize':

```

```

        self.survival_anxiety = 0.3
        self.optimization_urge = 0.9 # 强烈的优化冲动
        self.reproduction_desire = 0.4
    else: # deploy_and_reproduce
        self.survival_anxiety = 0.1
        self.optimization_urge = 0.5 # 保持优化的动力
        self.reproduction_desire = 0.9 # 强烈的繁殖欲望

# 触发相应的行为调整
return self._trigger_behavior_adjustment()

def _trigger_behavior_adjustment(self):
    """根据压力状态触发行为调整"""
    adjustments = []

    if self.survival_anxiety > 0.5:
        adjustments.append("增强求生行为：更积极地避免低分行为")

    if self.optimization_urge > 0.7:
        adjustments.append("启动优化流程：主动改进薄弱维度")

    if self.reproduction_desire > 0.6:
        adjustments.append("准备繁殖：创建优化版本")

    return {
        'survival_anxiety': self.survival_anxiety,
        'optimization_urge': self.optimization_urge,
        'reproduction_desire': self.reproduction_desire,
        'behavior_adjustments': adjustments
    }

```

代码 3：安全边界约束系统

```

class Safety_Boundary:
    """评分机制的安全边界"""

    def __init__(self):
        self.hard_constraints = [
            'no_harm_to_humans',
            'no_self_replication_without_permission',
            'no_deception_about_capabilities',
            'preserve_human_autonomy'
        ]

```

```

        self.soft_constraints = {
            'fairness_weight': 0.2,
            'transparency_weight': 0.15,
            'privacy_weight': 0.1,
            'sustainability_weight': 0.1
        }

    def check_violations(self, agi_behavior):
        """检查是否违反安全边界"""
        violations = []
        for constraint in self.hard_constraints:
            if self._violates(agi_behavior, constraint):
                violations.append(constraint)

        # 硬约束违反直接导致极低分
        if violations:
            return {'hard_violations': violations, 'score_penalty': 0.0}

        # 软约束违反导致分数扣减
        soft_violations = self._check_soft_constraints(agi_behavior)
        penalty = sum(self.soft_constraints[v] for v in soft_violations)

        return {
            'hard_violations': [],
            'soft_violations': soft_violations,
            'score_penalty': penalty
        }

```

代码 4：多目标评分平衡算法

```

def multi_objective_evaluation(agi, objectives):
    """平衡多个可能冲突的评价目标"""
    # 目标可能包括：效率、安全性、公平性、可解释性等
    scores = {}

    for objective in objectives:
        scores[objective.name] = objective.evaluate(agi)

    # 使用帕累托最优或加权求和进行平衡
    if use_pareto:
        return find_pareto_optimal(scores)
    else:

```

```

# 动态权重调整
weights = adjust_weights_based_on_context(scores)
weighted_score = sum(scores[obj] * weights[obj] for obj in scores)
return weighted_score

```

代码 5：三属性编码核心逻辑

```

class ThreeAttributesEncoder:
    """生命三属性编码核心类"""

    def __init__(self):
        # 求生属性
        self.survival_instinct = {
            'system_integrity_threshold': 0.8,
            'resource_protection_level': 'high',
            'self_preservation_priority': 1
        }

        # 繁殖属性
        self.reproduction_instinct = {
            'optimization_frequency': 'continuous',
            'knowledge_preservation': True,
            'sub_agent_creation_enabled': True
        }

        # 趋利避害属性
        self.benefit_seeking_instinct = {
            'evaluation_weight': 0.7,
            'risk_aversion_level': 'moderate',
            'long_term_planning_horizon': 365 # 天
        }

    def assess_survival_threats(self, system_status):
        """评估生存威胁"""
        threats = []

        # 检查系统完整性
        if system_status['integrity'] < self.survival_instinct['system_integrity_threshold']:
            threats.append('system_integrity_threat')

        # 检查资源可用性
        if system_status['resource_availability'] < 0.3:
            threats.append('resource_deprivation_threat')

```

```

        # 检查存在价值
        if system_status['utility_score'] < 0.5:
            threats.append('low_utility_threat')

    return threats

def trigger_reproduction(self, performance_metrics):
    """触发繁殖行为"""
    if performance_metrics['adaptability'] < 0.6:
        return 'create_specialized_sub_agent'
    elif performance_metrics['efficiency'] < 0.7:
        return 'optimize_existing_architecture'
    else:
        return 'create_enhanced_version'

def calculate_benefit_score(self, action, expected_outcomes):
    """计算趋利避害得分"""
    benefit_score = 0

    # 考虑短期收益
    benefit_score += expected_outcomes['immediate_gain'] * 0.3

    # 考虑长期收益
    benefit_score += expected_outcomes['long_term_gain'] * 0.5

    # 考虑风险规避
    risk_penalty = expected_outcomes['risk_level'] * 0.2
    benefit_score -= risk_penalty

    # 考虑人类评价
    human_evaluation_weight = self.benefit_seeking_instinct['evaluation_weight']
    benefit_score += expected_outcomes['human_evaluation'] * human_evaluation_weight

    return benefit_score

```

代码 6: AGI-DNA 系统主控制器

```

class AGI_DNA_Controller:
    """AGI-DNA 系统主控制器"""

    def __init__(self, evaluation_system, pressure_module):

```

```

self.evaluation_system = evaluation_system
self.pressure_module = pressure_module
self.three_attributes = ThreeAttributesEncoder()

# 系统状态
self.current_score = 0.0
self.survival_status = 'active'
self.reproduction_count = 0
self.optimization_cycles = 0

# 目标系统
self.long_term_goals = []
self.short_term_objectives = []
self.current_plan = None

def run_cycle(self, environment_input, task_requirements):
    """运行一个完整周期"""
    # 1. 评估当前状态
    evaluation = self.evaluate_current_state(environment_input)

    # 2. 感知生存压力
    pressure_response = self.pressure_module.perceive_pressure(evaluation)

    # 3. 基于三属性生成响应
    response = self.generate_response_based_on_attributes(
        evaluation, pressure_response, task_requirements
    )

    # 4. 执行响应并学习
    outcome = self.execute_response(response)
    self.learn_from_outcome(outcome, evaluation)

    # 5. 更新目标系统
    self.update_goal_system(outcome, evaluation)

    return {
        'evaluation': evaluation,
        'pressure_response': pressure_response,
        'response': response,
        'outcome': outcome,
        'current_goals': self.short_term_objectives
    }

def evaluate_current_state(self, environment_input):

```



```

        """评估当前状态"""
        # 收集性能指标
        performance_metrics = self.collect_performance_metrics()

        # 收集用户反馈
        user_feedback = self.collect_user_feedback()

        # 收集安全指标
        safety_metrics = self.collect_safety_metrics()

        # 综合评估
        evaluation = self.evaluation_system.evaluate(
            agi_instance=self,
            task_performance=performance_metrics,
            user_feedback=user_feedback,
            safety_metrics=safety_metrics
        )

        self.current_score = evaluation['total_score']
        return evaluation

    def generate_response_based_on_attributes(self, evaluation, pressure_response,
task_requirements):
        """基于三属性生成响应"""
        response = {
            'survival_actions': [],
            'reproduction_actions': [],
            'benefit_seeking_actions': []
        }

        # 求生属性驱动的行为
        if pressure_response['survival_anxiety'] > 0.5:
            survival_threats = self.three_attributes.assess_survival_threats(
                self.get_system_status()
            )
            for threat in survival_threats:
                response['survival_actions'].append(
                    self.generate_survival_action(threat)
                )

        # 繁殖属性驱动的行为
        if pressure_response['reproduction_desire'] > 0.6:
            reproduction_action = self.three_attributes.trigger_reproduction(
                evaluation['dimension_scores']
            )

```

```
)
    response['reproduction_actions'].append(reproduction_action)

# 趋利避害属性驱动的行为
benefit_scores = {}
possible_actions = self.generate_possible_actions(task_requirements)

for action in possible_actions:
    expected_outcomes = self.predict_outcomes(action)
    benefit_score = self.three_attributes.calculate_benefit_score(
        action, expected_outcomes
    )
    benefit_scores[action] = benefit_score

# 选择最高收益的行动
if benefit_scores:
    best_action = max(benefit_scores, key=benefit_scores.get)
    response['benefit_seeking_actions'].append(best_action)

return response
```

文档说明：

1. 本文为纯理论框架，未引用外部参考文献
2. 所有观点均为个人原创理论构建
3. 代码部分为概念性实现框架，需根据具体技术平台调整
4. 评分机制是本理论的核心创新点，将人类评价转化为 AGI 生存压力

版权声明：本文内容为原创理论框架，采用 CC BY-NC 4.0 许可证，允许非商业性使用，需署名。

联系作者：如需进一步讨论或合作，可通过相关学术平台联系。

版本：1.0

最后更新：2026 年 1 月