

# Contents

<b>AI 时代的人类能力建构与认知韧性理论</b>	<b>5</b>
摘要	5
Abstract	5
第一章 引言与理论定位	6
1.1 核心命题：评估 AI 的新标准	6
1.1.5 术语澄清与理论命名	6
1.2 理论缺口：现有范式的共同盲点	6
1.3 本理论的核心贡献：从诊断到解决方案	7
1.4 研究方法与论文结构	8
2.1 认知卸载研究：从现象描述到机制理解	9
2.1.1 认知卸载的基本概念	9
2.1.2 认知卸载的领域特异性证据	9
2.1.3 AI 时代的认知卸载：新兴证据	9
2.1.4 认知卸载的调节因子	10
2.1.5 认知卸载研究的理论空白	10
2.2 延展心智理论的批判性回顾：从哲学隐喻到可操作标准	10
2.2.1 延展心智理论的核心主张	10
2.2.2 延展心智理论的局限性	11
2.2.3 对等原则的重构：从功能对等到过程对等	11
2.3 自动化研究与脚手架理论：从失败案例到成功路径	11
2.3.1 自动化悖论：永久支持的代价	11
2.3.2 脚手架理论：渐进独立的教育智慧	12
2.3.3 两种模式的对比：外骨骼 vs 脚手架	12
2.4 神经可塑性与认知训练：适度挑战的生物学基础	12
2.4.1 用进废退：神经可塑性的双向性	12
2.4.2 最优挑战区：适度困难的学习曲线	13
2.4.3 认知训练争议与 CET 的澄清	13
2.5 理论整合：跨学科证据的收敛	13
2.5.1 证据链的收敛性：殊途同归的科学共识	14
2.5.2 CET 的独特贡献：从分散洞察到统一框架	14
2.5.3 承上启下：从“为什么到”“是什么”与如何做	15
<b>第三章 CET 核心理论构建</b>	<b>15</b>
3.0 核心术语与锚点定义	15
3.0.1 核心符号系统	15
3.0.2 反脆弱性验证原则 (AVP)【锚点 B1】	16
3.0.3 内共生最小法则 (EML)【锚点 B2】	16
3.0.4 LSA 功能分层【锚点 B3】	16
3.0.5 最优挑战区【锚点 B4】	16
3.0.6 边界条件【锚点 B5】	16
3.1 AVP 原则的深度阐释：从抽象判据到可操作测量	16
3.1.1 理论基础：从 Taleb 的反脆弱性到认知增强评估	16
3.1.2 AVP 的三级分级体系：从基础到卓越	17
3.1.3 AVP 测量协议：可复现的操作指南	18
3.1.4 跨领域校准：参数的情境敏感性	19
3.2 有益认知摩擦：挑战与支持的最优平衡	19
3.2.1 为什么需要认知摩擦：从神经科学到学习理论	19
3.2.2 最优挑战区：50-70% 成功率的量化目标	19
3.2.3 摩擦的动态调节：自适应算法的概念模型	20
3.3 系统性支持削减：从脚手架到能力独立	20
3.3.1 为什么必须削减：从脚手架理论到依赖预防	20
3.3.2 支持削减曲线：从 S4 到 S1 到 0 的演化路径	21

3.3.3 回退与加速机制: 保底安全线与提速路径	21
3.4 伙伴式主体性:AI 角色的哲学重构	22
3.4.1 从工具到伙伴:AI 角色的三种范式	22
3.4.2 伙伴式主体性的四个维度	22
3.4.3 伙伴关系的伦理边界	23
3.4.4 伙伴式主体性与其他 AI 范式的对比	23
3.5 认知外骨骼的病理学: 失败模式的系统分析	24
3.5.1 外骨骼的核心特征: 从健康到病理的转折点	24
3.5.2 外骨骼的 10 个预警信号	24
3.5.3 分级干预策略: 从预防到治疗	25
3.6 跨尺度概览: 从个体到团队到组织到社会	25
3.6.1 团队层面: 协作中的能力分布	25
3.6.2 组织层面: 认知基础设施与系统韧性	25
3.6.3 社会层面: 认知公地悲剧	26
3.6.4 跨尺度的同构性	26
第四章 跨尺度机制分析: 从个体到组织到社会	26
核心术语与符号表 (第四章)	26
4.1 个体 → 团队: 能力分布与协作模式的重构	27
4.1.1 从 I-AVP 到 T-AVP: 团队层面的反脆弱性	27
4.1.2 T-AVP 的测量协议	28
4.1.3 团队层面的设计启示	29
4.2 团队 → 组织: 认知基础设施的系统性风险	29
4.2.1 组织层面的新涌现: 制度化依赖	29
4.2.2 O-AVP: 组织反脆弱性验证	30
4.2.3 组织设计的三个关键机制	32
4.3 组织 → 社会: 认知公地悲剧	32
4.3.1 认知公地的概念: 集体行动困境	32
4.3.2 代际能力鸿沟: 三代人的分化	33
4.3.3 S-AVP: 社会认知资本的验证 (概念框架)	33
4.3.4 认知公地的潜在干预 (方向性讨论)	34
4.4 跨尺度共性机制: CET 的尺度不变性	34
4.4.1 AVP 原则的尺度不变性	34
4.4.2 级联脆弱: 从微观到宏观的传播	35
4.4.3 设计启示: 多尺度协同的必要性	36
4.5 小结: 理论的系统性与实践的紧迫性	37
4.5.1 本章核心贡献	37
4.5.2 与前后章的连接	37
4.5.3 理论的紧迫性	37
第五章 分层共生架构 (LSA): CET 的工程化实现	37
核心术语与符号表 (第五章)	38
5.1 LSA 总览: 为什么需要分层架构	38
5.1.1 传统 AI 系统的设计困境	38
5.1.2 LSA 的设计哲学: 能力建构优先	39
5.1.3 LSA 四层架构总览	39
5.2 L1 层: 基础 AI 能力层	41
5.2.1 职责与边界	41
5.2.2 L1 的典型能力模块	41
5.2.3 L1 与传统 AI 系统的区别	41
5.3 L2 层: 摩擦与削减层	41
5.3.1 核心挑战: 如何实现适度帮助	41
5.3.2 认知摩擦引擎 (CFE)	42
5.3.3 支持削减调度器 (SGS)	43

5.3.4 L2 层的工程挑战	45
5.4 L3 层: 监测与反馈层	45
5.4.1 核心职责: 能力追踪与 AVP 验证	45
5.4.2 AVP 遥测模块 (AVP-TM)	45
5.4.3 能力建模引擎	46
5.4.4 预警与干预系统	47
5.4.5 L3 层的隐私保护设计	47
5.5 L4 层: 编排与治理层	48
5.5.1 核心职责: 多尺度策略与伦理治理	48
5.5.2 多尺度编排器 (MSO)	48
5.5.3 伦理治理框架	50
5.5.4 全局优化目标函数	51
5.6 端到端工作流示例	52
5.6.1 个体层的 LSA 应用 (概念演示)	52
5.6.2 团队层的 LSA 应用 (概念演示)	53
5.7 技术可行性与工程挑战	54
5.7.1 现有技术栈的适配性	54
5.7.2 最小可行原型 (MVP) 设计方向	55
5.7.3 关键工程挑战	55
5.7.4 常见问题解答 (FAQ)	57
5.8 小结: LSA 的理论贡献与实践路径	58
5.8.1 本章核心贡献	58
5.8.2 与前后章的连接	59
5.8.3 LSA 的三个开放问题	59
5.8.4 实践路径建议	60
<b>第六章 局限性、证伪路径与未来方向</b>	<b>60</b>
6.1 理论的六大局限性	60
6.1.1 尺度边界: 个体以下与社会以上	60
6.1.2 任务类型的限制	61
6.1.3 测量挑战: 从概念到操作的鸿沟	62
6.1.4 个体差异与公平性悖论	63
6.1.5 技术依赖的双刃剑	65
6.1.6 文化嵌入性与价值判断	66
6.2 八个可证伪假说及其证伪路径	67
6.2.1 核心假说概览与因果识别策略	67
6.2.2 个体层假说的证伪路径	69
6.2.3 团队与组织层假说的证伪路径	70
6.2.4 LSA 设计假说的证伪路径	72
6.2.5 社会层假说的证伪路径	72
6.3 未来研究议程: 三个时间尺度	73
6.3.1 短期研究 (1-3 年): 验证核心机制	73
6.3.2 中期研究 (3-5 年): 跨尺度扩展与理论整合	74
6.3.3 长期研究 (5-10 年以上): 社会影响与理论演化	76
6.4 研究伦理与开放科学承诺	77
6.4.1 伦理原则	77
6.4.2 开放科学承诺	77
6.5 结语: 理论的生命在于批判与演化	78
<b>第 7 章 - 术语与符号系统</b>	<b>78</b>
1. 核心概念固定锚点 (B1-B5)	78
B1   AVP 定义锚点	78
B2   EML 定义锚点	79
B3   LSA-F 功能分层锚点	79

B4   最优挑战区锚点	79
B5   边界条件锚点	79
2. 核心参数速查表	79
3. 缩写速查表	80
4. 核心评估原则	80
4.1 等效努力原则 (Equivalent Effort Principle)	80
4.2 Goodhart 防护原则	81
5. 核心术语中英对照	81
5.1 理论核心概念	81
5.2 测量相关术语	82
5.3 架构相关术语	82
6. 参数登记簿 (单一事实源 - Single Source of Truth)	82
7. 文档维护协议	83
7.1 更新流程	83
7.2 一致性检查命令	83
7.3 版本历史	83
<b>CET8 - 图表库与对照表</b>	<b>83</b>
图表库使用说明	83
B.1 外骨骼 vs 内共生核心对照表	84
B.2 跨尺度验证框架 ( $I \rightarrow T \rightarrow O \rightarrow S$ )	84
B.3 支持削减曲线对比	85
B.4 最小 AVP 实验流程	86
B.5 能力向量 $C(t)$ 概念模型	87
B.6 LSA 四层架构示意	87
B.7 图表素材技术规范 (便于复用与派生)	88
B.8 图表设计的可访问性自检清单	89
<b>CET9 - 综合附录系统</b>	<b>89</b>
附录 A: AVP 测量协议工具包	89
A.0 工具包定位与声明	89
A.1 标准化 AVP 测量清单	89
A.5 常见问题与注意事项	93
A.6 简化版协议: 24 小时拔线轻量测试	95
附录 B: 补充案例分析	95
B.0 案例选择方法论	95
B.1 外骨骼案例: 写作助手的依赖陷阱	96
B.2 内共生案例: 编程教学平台的成功实践	97
B.3 团队层案例: 软件公司的 T-AVP 实验	99
附录 C: 理论定位与学术对话	100
C.1 CET 与认知心理学的对话	100
C.2 CET 与教育技术的对话	101
C.3 CET 与 AI 伦理的对话	101
C.4 CET 与组织行为学的对话	102
C.5 CET 与系统科学的对话	102
附录 D: 术语与符号索引	103
D.1 核心概念术语 (按字母序)	103
D.2 符号与参数索引	103
D.3 缩写速查表 (完整版)	104
D.4 跨章节概念地图	105
参考文献	105

# AI 时代的人类能力建构与认知韧性理论

*Adjust task format without lowering challenge intensity; assess by relative improvement rather than absolute level; (if accessibility applies) challenge budget conservation.*

副标题：基于反脆弱性验证原则的认知内共生范式

作者：杨国平

邮箱：a44425874@gmail.com

版本：V1.0

日期：2025 年 9 月

## 摘要

研究目的：随着生成式 AI 的广泛应用人类独立认知能力面临退化风险本研究旨在建立一个可证伪的理论框架用于评估和优化 AI 时代的人类能力建构

研究方法：核心包括三大机制：(1) 反脆弱性验证原则（AVP）——通过“拔线测试检验协作是否促进独立能力”(2) 内共生最小法则（EML）——设计有益认知摩擦与系统性支持削减(3) 伙伴式主体性——AI 作为认知伙伴的角色定位采用跨学科综合方法整合认知心理学、教育技术、组织行为学和 AI 伦理的研究成果并通过文献分析、案例研究和概念建模进行验证

主要结果：建立了四层跨尺度验证体系识别了“认知外骨骼（过度依赖导致能力退化）与“认知内共生（能力持续增强）两种模式并设计了分层共生架构（LSA）作为工程化实现路径通过三个领域的案例分析（写作、编程、团队协作）验证了理论的解释力和预测力研究明确提出六大局限性和八个可证伪假说（H1-H8）为后续实证研究提供清晰路径

研究结论：为 AI 时代的能力建构提供了首个系统性、可操作的评估标准其价值在于：提出了超越表面效率、关注独立能力的评估范式建立了跨尺度的一致性框架设计了可工程化的技术实现路径理论的所有参数均为工作假设需通过跨领域实证研究校准我们主动承认理论局限并期待学术社区的批判、验证和超越

关键词：人工智能认知能力建构反脆弱性验证人机协作能力评估认知韧性

## Abstract

**Objective:** With the widespread adoption of generative AI human independent cognitive capabilities face risks of degradation. This study aims to establish a falsifiable theoretical framework for assessing and optimizing human capability building in the AI era.

**Methods:** The core comprises three mechanisms: (1) Antifragility Validation Principle (AVP)—verifying whether collaboration enhances independent capability through “unplugged tests”; (2) Endosymbiotic Minimal Law (EML)—designing beneficial cognitive friction and systematic support reduction; (3) Partner-like Agency—positioning AI as a cognitive partner. The study adopts an interdisciplinary synthesis approach integrating research from cognitive psychology educational technology organizational behavior and AI ethics validated through literature analysis case studies and conceptual modeling.

**Results:** This paper establishes a four-layer cross-scale validation system identifies two contrasting patterns —“cognitive exoskeleton” (over-reliance leading to capability degradation) versus “cognitive endosymbiosis” (continuous capability enhancement)—and designs a Layered Symbiosis Architecture (LSA) as an engineering implementation pathway. Through case analyses across three domains (writing programming team collaboration) the theorys explanatory and predictive power is validated. The study explicitly articulates six major limitations and eight falsifiable hypotheses (H1-H8) providing clear pathways for subsequent empirical research.

**Conclusions:** This theory provides the first systematic operational assessment standard for capability building in the AI era. Its value lies in: proposing an assessment paradigm that transcends surface efficiency to focus on independent capability; establishing a consistent framework across scales; designing an engineerable

technical implementation pathway. All theory parameters are working hypotheses requiring cross-domain empirical calibration. We actively acknowledge theoretical limitations and welcome critique validation and transcendence from the academic community.

**Keywords:** Artificial Intelligence; Cognitive Capability Building; Antifragility Validation; Human-AI Collaboration; Capability Assessment; Cognitive Resilience

## 第一章 引言与理论定位

### 1.1 核心命题：评估 AI 的新标准

当我们评估一个 AI 工具的价值时我们真正应该关注什么

主流的评价标准聚焦于使用时的表现：任务完成速度、输出质量、用户满意度然而本文提出一个根本性的反驳：**AI** 工具的真正价值不在于使用时你有多强而在于离开它时你有多强

我们将这一原则形式化为反脆弱性验证原则（Antifragility Validation Principle AVP）：

以拔线测试（Unplugged Test）检验协作是否促进独立能力核心判据： $P_2 > B_0 +$

详细定义见 3.0.2 节（AVP 定义锚点 B1）

一句话原则：设计靠 **EML**（摩擦 + 削减）| 验收靠 **AVP**（拔线 + 对比）

这个判据揭示了两种根本不同的 AI 使用范式：

场景对比：想象一位医生、程序员或学生在使用 AI 辅助工作 6 个月后突然失去 AI 访问权限：

- 认知外骨骼模式：其独立表现显著下降（ $P_2 < B_0$ ）甚至无法完成原本能够胜任的任务
- 认知内共生模式：其独立表现不仅保持甚至超越原有水平（ $P_2 > B_0 +$ ）因为 AI 的使用过程强化了其底层能力

这不是假设场景神经科学研究已经提供了警示信号：Dahmani 等人（2020）的研究发现习惯性使用 GPS 导航的人群其空间记忆任务表现与海马体灰质体积存在显著关联在教育领域多项初步研究共同指向一个趋势：AI 写作工具的重度用户在脱离辅助时其独立写作的流畅度与结构组织能力呈现下降趋势（具体效应量因研究设计差异而变化但方向具有一致性）这种现象在认知心理学中被称为认知卸载 “(cognitive offloading)” 其负面效应已有充分的实证基础（Risko & Gilbert 2016; Sparrow et al. 2011）

更令人警醒的是这种退化的隐蔽性每个用户都感觉自己在进步——任务完成更快输出质量更高但这种表面的能力提升可能掩盖了底层认知架构的系统性侵蚀当我们将视角拉长到代际尺度一个严峻的问题浮现：如果主流 **AI** 使用模式持续遵循外骨骼逻辑人类社会可能面临一种新型的文明脆弱性——不是因为 **AI** 背叛我们而是因为我们主动放弃了独立思考的能力

#### 1.1.5 术语澄清与理论命名

关于“认知外骨骼与认知内共生”：

本文使用“认知外骨骼”指代 AI 使用导致人类能力退化的病理模式“认知内共生指代 AI 使用促进人类能力提升的健康模式 理论的核心目标是避免外骨骼、建立内共生

术语层级：- **AVP**：评估是否达成内共生的验收标准 - **EML**：设计内共生系统的构造法则 - 内共生：健康的人机关系（理论目标）- 外骨骼：病理的依赖状态（理论警示）

### 1.2 理论缺口：现有范式的共同盲点

当前的人机交互研究和实践主要由三种范式主导但它们都未能有效应对上述挑战：

**1.2.1 工具范式（Tool Paradigm）** 将 AI 视为被动的效率提升工具强调人类的完全主导地位这一范式（如传统 HCI 研究）忽视了长期使用对用户认知模式的反向塑造作用假设工具是中性的使用者始终保持独立性

**1.2.2 增强范式 (Augmentation Paradigm)** 如 Engelbart (1962) 的智力增强愿景关注如何通过技术放大人类能力但这一范式往往假设“增强”必然是正面的缺乏对“负向”认知增益效应的识别机制

Extended Mind 理论 (Clark & Chalmers 1998) 虽然承认认知的延展性但未能区分良性延展 (促进成长) 与病理延展 (导致依赖) 本理论的贡献之一正是提供了可证伪的阈值来区分这两种延展模式

**1.2.3 自动化范式 (Automation Paradigm)** 聚焦用 AI 替代人类完成任务以实现效率最大化这一范式对人类长期认知健康几乎不予考虑“自动化悖论” (Parasuraman & Riley 1997) 和“自动化自满 (automation complacency) 研究已经警示了过度依赖自动化系统的风险但这些洞察尚未系统地应用于 AI 工具设计

这三种范式的共同盲点在于：它们都将 AI 视为外在于人类认知系统的工具或替代品而未能将人机关系置于一个动态、共生、相互塑造的系统框架中理解更关键的是现有研究缺乏一个可操作的、可证伪的标准来评估人机交互的长期健康性我们如何判断一个 AI 系统是在赋能用户还是在制造依赖这一核心问题在学术和工业界都缺乏清晰答案

### 1.3 本理论的核心贡献：从诊断到解决方案

本研究旨在填补上述理论缺口为 AI 时代的人机协作提供一个可验证的范式框架理论贡献体现在四个层面：

**1.3.1 建立判别标准：反脆弱性验证原则 (AVP)** 本研究的首要贡献是将 AVP 确立为评估人机交互健康性的参考标准这不仅是一个技术指标更是一个价值判断的范式革命：它挑战了科技行业长期以即时效率为核心的评价体系

AVP 的可操作化包括：

- 基线测量 ( $T_0$ )：用户使用 AI 前的独立能力
- 协作测量 ( $T_1$ )：用户与 AI 协作时的表现
- 独立测量 ( $T_2$ )：用户在拔线窗口 ( $W=4-8$  周默认 6 周工作假设需跨领域/任务校准) 后的独立能力

测量要求： $T_0$  与  $T$  应采用等值平行测验以控制重测与熟悉化效应

判定标准：

- $P_2 \geq B_0 +$ ：成功 (认知内共生)
- $P_2 \approx B_0$ ：中性 (未造成损害但也未促进成长)
- $P_2 < B_0$ ：失败 (认知外骨骼造成依赖性退化)

**1.3.2 确立设计原则：内共生最小法则 (EML)** 本研究提出内共生最小法则 (Endosymbiotic Minimal Law) 作为健康人机协作的判定框架：

内共生最小法则 (EML)：构成“认知内共生的设计必要条件为：

- (1) 有益认知摩擦：使用户处于最优挑战区 (成功率 50-70% 工作假设需跨领域/任务校准个体自适应)
- (2) 系统性支持削减：AI 支持强度按既定削减曲线从  $S_4 \rightarrow S_1 \rightarrow S_0$

二者为联合充分的设计条件但最终仍需 AVP ( $P_2 \geq B_0 +$ ) 作为验收必要条件

逻辑关系：条件 (1) 和 (2) 构成设计层面的联合充分条件 AVP 验证是结果验收的必要条件三者联合构成内共生的充分必要条件任一条件缺失系统即判定为认知外骨骼

边界条件声明：

边界条件：本理论适用于能力增强型人机协作补偿性外骨骼 (如残障辅助、超越生理极限的设备) 不适用此判据所有参数均为概念工作模型需跨领域校准

**1.3.3 重新定义 AI 角色：从工具到认知伙伴** 本研究提出将 AI 从被动工具重塑为具有伙伴式主体性（Partnership Agency）的认知共生体这种转变的核心不是拟人化而是功能性的角色重构：

可操作化的三个锚点：

1. 摩擦注入：AI 主动创造适度认知挑战（而非总是提供最简单路径）
2. 脚手架消退：遵循系统性支持削减曲线
3. 以 AVP 闭环为交互终点：协作的最终目标是用户独立能力的提升

这种“伙伴性与传统 HCI 追求的无缝、零摩擦”体验形成根本性对立：一个真正的认知伙伴不应只是顺从的助手更应是能够提出建设性挑战的协作者

**1.3.4 技术实现框架：分层共生架构（LSA）** 本研究提出分层共生架构（Layered Symbiosis Architecture LSA）作为技术承载第一章仅概述核心结构详细实现见第五章：

**LSA-F（功能分层）：** L1 知识整合 | L2 状态建模 | L3 摩擦校准 | L4 元认知协调

支持档位栈（ $S_4 \rightarrow S_1 \rightarrow S_0$ ）用于表达支持强度与 LSA-F 为正交维度

表 1：认知外骨骼 vs 认知内共生的核心对照

维度	认知外骨骼	认知内共生
设计哲学	替代/卸载	赋能/强化
认知摩擦	最小化	优化（50-70% 工作假设）
时间性	永久依赖	临时共生
支持削减	无/固定支持	系统性递减（ $S_4 \rightarrow S_1 \rightarrow S_0$ ）
AVP 结果	$P_2 \quad B_0$	$P_2 \quad B_0 +$
神经效应 *	能力退化趋势	能力增强趋势

\* 注：神经效应一栏为基于现有认知神经科学研究（如 Dahmani et al. 2020; Maguire et al. 2000）的趋势性推断需要进一步针对性实证研究验证

**1.3.5 理论定位：规范性解决方案框架** 本研究定位为一个规范性解决方案框架（normative solution framework）我们不仅诊断问题更明确应当如何构造人机协作以避免依赖退化并用 AVP 作为参考标准验证是否达标

与 AI 对齐研究的关系：

本理论与 AI 对齐（AI Alignment）研究是互补而非对立的关系：

- 对齐研究关注 AI 的意图是否与人类价值一致
- 本理论关注人机协作是否促进人类能力的可持续发展

我们提供的是在假设 AI 已对齐的前提下如何设计健康的人机交互模式的完整方案包括：

- 构造标准（通过 EML）
- 验收标准（通过 AVP）
- 工程路径（通过 LSA）

## 1.4 研究方法与论文结构

**1.4.1 方法论定位** 本研究采用理论构建与概念分析相结合的方法本理论定位为一个可证伪的理论框架：我们提出了一套可操作的概念体系和八个可测试的假说（见第六章）但坦诚承认当前阶段的理论尚未经过大规模实证验证

透明性声明：

1. 本文使用的所有量化参数（如 Cohens d 0.3 或 10%（*working assumption*）成功率 50-70% 拔线窗口 W=4-8 周默认 6 周）均为概念工作模型基于认知心理学和教育测量学文献的合理推断但需要跨领域校准和实证检验



2. 案例选择遵循理论启发性标准不追求统计代表性而是用于阐释机制和边界条件
3. 我们明确指出理论的适用边界和证伪路径（详见第六章）坚持开放科学原则

**1.4.2 论文结构** 本文共六章：第二章梳理跨学科证据基础第三章构建 AVP/EML 理论框架 第四章扩展至团队和组织层面第五章讨论技术实现路径 第六章明确理论边界并提出证伪路径 # 第二章 文献综述与理论基础

本章旨在系统梳理 CET 理论的多学科基础揭示现有研究的拼图如何指向一个统一的理论框架我们将展示：虽然认知卸载、自动化悖论、神经可塑性等领域各自提供了重要洞察但它们缺乏一个可操作的评估标准和统一的设计原则——这正是 CET 理论的核心贡献所在

## 2.1 认知卸载研究：从现象描述到机制理解

### 2.1.1 认知卸载的基本概念

认知卸载（Cognitive Offloading）指的是个体将认知任务委托给外部系统（如工具、技术或环境）以减轻内部认知负荷的现象 Risko 和 Gilbert（2016）在其权威综述中将认知卸载定义为使用物理行动来改变信息处理需求的策略这一概念整合了多个相关研究领域包括分布式认知、延展心智和具身认知

早期关于认知卸载的开创性研究来自 Sparrow 等人（2011）对 Google 效应的探索他们通过一系列实验发现当参与者知道信息可以在计算机上轻松获取时他们对信息本身的记忆显著下降但对信息存储位置的记忆有所增强这项研究首次系统性地揭示了数字技术对记忆模式的重塑作用引发了关于互联网是否改变了我们记忆方式的广泛讨论

Barr 等人（2015）进一步推进了这一研究提出智能手机被用作“认知替代品（Cognitive Substitutes）”他们的研究表明认知卸载不仅影响记忆更广泛地影响思维过程本身：当面对需要深度思考的问题时重度智能手机使用者更倾向于立即求助于搜索引擎而非进行内部推理这种模式与分析性思维能力呈负相关（工作假设基于横断研究需纵向验证）

### 2.1.2 认知卸载的领域特异性证据

空间认知领域:**GPS** 导航的案例

GPS 导航对空间认知的影响提供了清晰的神经生物学证据 Dahmani 和 Bohbot（2020）的研究发现习惯性使用 GPS 导航的个体在自主导航任务中表现出空间记忆缺陷这种缺陷与海马体灰质体积存在关联（方向性证据需控制自选偏差）

这与 Maguire 等人（2000）对伦敦出租车司机的研究形成鲜明对比：出租车司机通过长期空间记忆训练其海马体后部显著增大这两项研究共同构成了一个“自然实验”：同一认知功能（空间导航）在不同的技术关系下（独立训练 vs GPS 依赖）导致相反的神经可塑性结果这为 CET 理论关于技术关系决定认知后果的核心主张提供了神经层面的支持

计算与数学领域

计算器的普及引发了关于心算能力的争议教育研究显示在计算器广泛使用后成长的学生群体其心算流畅度和数感呈现变化趋势（方向性证据因测量标准、教学方法和对照组设置而异结论尚不一致）

### 2.1.3 AI 时代的认知卸载：新兴证据

近年多项研究与预印本在 AI 辅助写作与信息检索场景中观测到“更频繁的技术性卸载”与更弱的独立批判性思维/结构化写作之间存在显著关联这些研究采用批判思维评估量表（如 HCTA）、写作质量评分、以及认知卸载行为自陈量表等多种测量工具（Noy & Zhang 2023; Peng et al. 2023; DellAcqua et al. 2023）

虽然具体效应量因研究设计（横断 vs 纵向）、样本特征和测量方式而异但方向总体一致：更高度度的认知卸载与更弱的独立认知表现相关联（工作假设需排除反向因果可能性）这些结果与 Google 效应等早期数字记忆研究相呼应提示卸载模式可能改变深加工策略的使用频率在教育场景中这种双刃剑效应尤为明显最新研究揭示了 AI 认知悖论现象：AI 工具既能增强学习效率也可能侵蚀批判性思维能力对 206 名职业教育学生的研究发现过度依赖 AI 辅助显著影响了学生的认知深度和问题解决能力这正是 CET 理论试图通过 AVP 测量和 EML 设计解决的核心问题

方法学注意事项: 鉴于现有证据大多为横断或短期观察存在自选偏差、反向因果（能力弱者更倾向卸载）等方法学威胁我们将其作为方向性经验依据而非定量结论理想的验证设计应采用纵向/**RCT** 并包含拔线窗口（**W=4-8** 周默认 **6** 周）的设计这正是 CET 理论主张的 AVP 验证路径（见第 3.3 节）

#### 2.1.4 认知卸载的调节因子

认知卸载并非一律有害研究表明其效应受多个因素调节:

任务阶段: 当任务处于新手阶段时适度的认知卸载可以降低工作记忆负荷促进初步理解但若在整个学习过程中持续依赖卸载而不逐步内化则可能导致能力发展停滞

卸载类型:AI 提供的是检索线索/提示（如关键词提醒）还是整段替代推理（如完整答案生成）其对能力发展的影响截然不同前者可能促进深度加工后者更可能绕过认知努力

过程设计: 训练方案是否包含系统性支持削减、是否要求回忆/迁移测验、是否提供延迟再测等因素决定了卸载是成为通往独立的脚手架还是演变为永久的“拐杖”

**CET** 的视角: 这些调节因子正对应 EML 的核心条件（见第 3.4 节）零摩擦的生成式替代在缺乏有益摩擦、系统性支持削减和 AVP 验证的情况下更可能在拔线后暴露能力缺口因此评估认知卸载的健康性不能只看卸载本身更要看其嵌入的交互设计是否符合 **EML** 原则

#### 2.1.5 认知卸载研究的理论空白

尽管认知卸载研究取得了丰富的实证成果但存在三个关键的理论空白:

空白 1: 缺乏评估标准

现有研究主要是描述性的: 它们揭示了认知卸载现象的存在、测量了其程度和相关因素但没有提供规范性的判断标准什么程度的认知卸载是有益的什么时候从合理利用工具转变为有害依赖 Risko 和 Gilbert（2016）的综述虽然全面但未能解决这一评估问题

**CET** 的贡献:AVP 原则（ $P_2 \ B_0 +$  见第 3.3 节）提供了第一个可操作的、可证伪的评估标准关键洞察在于: 不能仅看  $P_1$ （使用 AI 时的表现）必须测量  $P_2$ （拔线后的独立表现）

空白 2: 缺乏设计指导

认知卸载研究告诉我们什么可能是问题但没有系统地回答如何设计解决方案简单的”节制使用”建议既不现实也不可行——在数字化世界中完全避免技术使用是不可能的

**CET** 的贡献:EML（内共生最小法则见第 3.4 节）提供了从问题诊断到解决方案设计的完整路径明确了有益 AI 使用的必要条件: 有益摩擦、系统性支持削减、AVP 验证

空白 3: 缺乏跨尺度整合

现有研究主要聚焦于个体层面的认知卸载对组织和社会层面的影响缺乏系统性探索即使讨论社会影响也往往停留在抽象的文化批判层面缺乏机制性的解释框架

**CET** 的贡献: 通过跨尺度的统一框架将个体认知卸载、组织依赖锁定置于同一套理论原理之下提供了从微观到宏观的连贯解释（详见第四章）

## 2.2 延展心智理论的批判性回顾: 从哲学隐喻到可操作标准

### 2.2.1 延展心智理论的核心主张

Clark 和 Chalmers（1998）在其开创性论文《延展心智》（The Extended Mind）中提出了一个颠覆性的主张: 认知的边界不必局限于颅骨或皮肤当外部工具以适当的方式与认知过程耦合时它们应被视为认知系统的一部分

经典案例是 Otto 的笔记本:Otto 患有阿尔茨海默症他依赖笔记本记录信息当他查阅笔记本寻找博物馆地址时这个过程在功能上等同于健康人从记忆中提取信息 Clark 和 Chalmers 的核心论证是对等原则（Parity Principle）:

如果在面对某个任务时世界的某个部分发挥的功能若在头脑中完成我们会毫不犹豫地认为它是认知过程的一部分那么该部分就是认知过程的一部分——即使它在头脑之外

这一理论引发了认知科学、哲学、人机交互领域的广泛讨论并衍生出延展认知（Extended Cognition）、“认知整合”（Cognitive Integration）等相关概念

### 2.2.2 延展心智理论的局限性

尽管延展心智理论具有启发性但它在应用于 AI 时代的人机协作时存在两个根本性局限：

局限 1：缺乏健康性判别标准

延展心智理论关注的是有什么可以算作认知但没有回答什么样的延展是健康的按照对等原则 Otto 的笔记本和一个完全依赖 AI 完成所有思考任务的人都可以被视为“延展认知”的例子但这两种情况的健康性显然不同：

- Otto 的笔记本：补偿性工具帮助他维持生活自理
- 过度依赖 AI：能力退化的路径导致独立思考能力萎缩

延展心智理论没有提供区分这两者的标准因为它是一个描述性理论而非规范性理论

**CET** 的补充：AVP 原则为“良性延展与病理延展”提供了可证伪的分界：如果延展导致  $P_2 > B_0 + \Delta$ （拔线后能力提升）则为良性如果  $P_2 < B_0 - \Delta$ （拔线后能力下降）则为病理

局限 2：忽视过程的时间性

延展心智理论关注的是某一时刻的认知状态但人类能力的发展是一个时间过程一个工具在某一时刻可能是有益的延展但长期使用后可能导致能力萎缩

例如：

- 阶段 1：AI 辅助写作帮助新手快速产出此时为有益延展
- 阶段 2：持续使用导致独立写作能力停滞此时已演变为依赖

延展心智理论无法捕捉这种动态转变因为它缺乏过程视角和能力发展维度

**CET** 的补充：EML 的系统性支持削减条件明确要求工具的支持强度应从  $S_4 \rightarrow S_1 \rightarrow S_0$  确保延展是临时的脚手架而非永久的拐杖最近的实证研究进一步验证了这一现象通过两阶段实验研究者开发了 7 种认知参与技术发现缺乏认知摩擦的 AI 辅助会导致表面学习和技能发展障碍学习者产生学习幻觉——误以为自己掌握了知识实际上只是依赖工具完成了任务（Liao et al. 2024）这一发现与 CET 理论的 H2 假说高度一致表明认知摩擦在 AI 辅助学习中的关键作用

### 2.2.3 对等原则的重构：从功能对等到过程对等

CET 对延展心智理论的改造不是否定其核心洞察而是将其从“功能对等”升级为过程对等：

延展心智（Clark & Chalmers）：关注工具在某一时刻是否发挥认知功能 ↓ **CET** 的重构：关注工具的使用过程是否促进长期能力发展

这种重构的关键在于：将对等原则的判断标准从静态功能转向动态过程不仅要问工具是否帮助完成任务更要问工具的使用是否促进用户能力提升

小结：延展心智理论为 CET 提供了哲学基础但 CET 通过引入 AVP/EML 为这一宏大框架增加了可证伪的操作阈值具体而言：延展心智提供定义（什么算作认知延展）**AVP** 提供健康性判别阈值（良性 vs 病理延展的可测标准）**EML** 提供过程约束（如何设计促进能力的延展系统）三者合成了可证伪、可设计、可工程化的人机共生框架

## 2.3 自动化研究与脚手架理论：从失败案例到成功路径

### 2.3.1 自动化悖论：永久支持的代价

自动化研究领域提供了丰富的负面证据：过度自动化如何导致操作员技能退化 Bainbridge（1983）的开创性论文《自动化的讽刺》（Ironies of Automation）指出：自动化越完善操作员在需要介入时的能力越差

航空案例：飞行员技能退化

现代民航飞机高度自动化飞行员在正常飞行中主要处于监督角色但当自动系统失效时飞行员需要手动操作的技能往往已经退化多起航空事故（如法航 447 航班）的调查显示（方向性证据基于事故调查报告与行业研究）：

- 飞行员在自动驾驶时表现正常 ( $P_1$  高)
- 系统失效后手动操作能力不足 ( $P_2$  低)
- 长期依赖自动化与应急处理能力萎缩高度相关

这种永久性支持 → 能力退化的模式在多个独立研究与事故调查中反复出现正是 CET 警示的认知外骨骼模式: 永久的支持导致独立能力退化在拔线时暴露脆弱性

### 2.3.2 脚手架理论: 渐进独立的教育智慧

与自动化悖论形成对比的是教育心理学的”脚手架理论”(Scaffolding Theory Wood et al. 1976) 该理论强调: 有效的教学支持应该是临时的、渐进撤出的

脚手架理论的核心原则:

1. 最初支持: 在学习者能力不足时提供密集支持
2. 渐进淡出 (Fading): 随着学习者能力提升逐步减少支持
3. 最终独立: 支持完全撤出学习者能独立完成任务

这一模式在教学实践中被广泛验证: 从幼儿识字 (家长逐步减少指读) 到职业培训 (师傅逐步放手) 都遵循这一原则 (方向性证据基于教育学经验需控制混杂因素)

**CET** 的整合: 脚手架理论的渐进淡出正是 EML 的系统性支持削减条件的理论来源但脚手架理论缺乏定量的评估标准 CET 通过 AVP 补充了这一空白: 淡出是否成功不看学习过程 ( $P_1$ ) 而看独立表现 ( $P_2$ )

### 2.3.3 两种模式的对比: 外骨骼 vs 脚手架

为便于理解 CET 如何整合自动化研究的负面教训与脚手架理论的正面经验表 2.1 按照同构维度并列展示自动化失败路径与”脚手架成功路径”; 二者在 CET 中分别对应”认知外骨骼”(病理模式) 与认知内共生 (健康模式) 详见 3.0 节 AVP/EML 核心定义)

表 2.1: 自动化与脚手架的范式对比

维度	自动化模式 (失败路径)	脚手架模式 (成功路径)	CET 术语
支持时长	永久性	临时性	认知外骨骼 vs 认知内共生
支持强度	固定或递增	渐进递减至 0	$S_4 \rightarrow S_1 \rightarrow S_0$ 削减 vs 固定支持
设计目标	替代人类操作	促进能力发展	效率优先 vs 能力优先
失败表现	$P_1$ 高但 $P_2$ 低	$P_1$ 适中但 $P_2 B_0 +$	外骨骼症状 vs AVP 通过
典型案例	飞行员技能退化	学徒制培训成功	见第 4.X 节案例

关键洞察: 自动化研究揭示了”永久支持 → 能力退化”的高度风险脚手架理论指出了渐进独立的成功路径 CET 将两者整合为统一框架:EML 的两个条件分别对应避免自动化陷阱 (有益摩擦) 和实现脚手架效果 (支持削减)

## 2.4 神经可塑性与认知训练: 适度挑战的生物学基础

### 2.4.1 用进废退: 神经可塑性的双向性

神经科学研究揭示了大脑的用进废退原则: 经常使用的神经通路被强化长期闲置的通路被削弱这一机制为 CET 理论提供了生物学基础

正向证据: 伦敦出租车司机

Maguire 等人 (2000) 的经典研究发现伦敦出租车司机经过 3-4 年的复杂空间记忆训练后海马体后部灰质体积显著增大这种变化与驾驶年限正相关且在退休后有所回退 (方向性证据基于横断和纵向数据)

负向证据:GPS 依赖者

Dahmani 和 Bohbot (2020) 的研究发现习惯性使用 GPS 导航的个体在空间记忆任务中表现较差且海马体灰质体积较小 (方向性证据需控制自选偏差)

**CET** 的解释: 这两项研究共同验证了 EML 的核心机制:

- 有益摩擦（独立导航训练）→ 神经强化（海马体增大）
- 零摩擦卸载（GPS 依赖）→ 神经退化（海马体减小）

#### 2.4.2 最优挑战区：适度困难的学习曲线

认知心理学的意欲的困难（Desirable Difficulties）理论（Bjork 1994）指出：适度的学习困难能促进长期保持和迁移这一理论有丰富的实证支持：

- 间隔效应（Spaced Practice）：分散练习优于集中练习
- 测试效应（Testing Effect）：检索练习优于重复阅读
- 交错效应（Interleaving）：混合练习优于分块练习

这些效应的共同特点是：短期表现可能较差（ $P_1$  低）但长期保持和迁移更好（ $P_2$  高）这与 CET 的核心主张完全一致：不能只看协作时的表现（ $P_1$ ）必须测量拔线后的表现（ $P_2$ ）

最优挑战区的量化：

认知负荷理论（Cognitive Load Theory Sweller 1988）和教育心理学研究提示学习任务的成功率 50-70%（工作假设需跨领域/任务校准）\*\* 区间时学习效果最优（工作假设需跨领域/任务校准基于多个独立研究的方向性证据）：

- >85% 成功率：任务过于简单接近认知卸载无足够挑战
- 50-70% 成功率：适度困难促进深度加工和模式学习
- <30% 成功率：任务过难导致挫败和回避行为

这一区间为 EML 的”有益认知摩擦条件提供了定量参考”（见第 3.4 节详细讨论）

#### 2.4.3 认知训练争议与 CET 的澄清

认知训练领域存在关于”远迁移”效果的持续争议一些研究报告认知训练（如工作记忆训练）能提升其他认知能力另一些研究则未能复现这些效果（Melby-Lervåg & Hulme 2013; Simons et al. 2016）

**CET** 对争议的重构：

1. 有摩擦是必要非充分条件：仅有认知挑战不足以保证能力迁移需要配合系统性支持削减和 AVP 验证
2. 测量焦点的转移：传统训练研究往往只测量训练期间的表现提升（类似  $P_1$ ）但未能系统地评估训练停止后在新情境下的独立表现（ $P_2$ ）
3. 远迁移的困难正说明：避免出现  $P_1$  高、 $P_2$  不增的表面学习假象需要更严格的验证——这正是 AVP 的价值所在

设计启示：从神经科学到工程参数

本节的神经科学和认知训练证据可以转化为三个可操作的设计旋钮：

1. 摩擦旋钮：通过自适应算法将任务难度维持在成功率 **50-70%** 的最优挑战区（工作假设需跨领域/任务校准个体需动态校准）
2. 削减旋钮：AI 支持强度从约 80% 线性或分段递减至 0% 节律根据用户表现动态调整
3. 验证闭环：使用等值平行测验在  $T_0$  和  $T_3$ （拔线窗口 **W=4-8** 周默认 **6** 周之后）时间点测量判据为  $P_2$   
 $B_0 +$  (Cohens d 0.3 或 10% (*working assumption*) 工作假设需跨领域/任务校准)

注：以上参数均为基于文献的工作假设需要跨任务类型、领域和人群进行实证校准验证

### 2.5 理论整合：跨学科证据的收敛

本章通过四个独立研究领域的系统梳理为 CET 理论构建了坚实的实证基础这些来自认知心理学、哲学、工程心理学、教育学和神经科学的证据虽然起源不同、方法各异但却惊人地收敛于相同的核心洞察：人与工具的关系质量决定了人类能力的演化方向本节将整合这些证据链明确 CET 如何填补现有研究的理论空白

2.5.1 证据链的收敛性：殊途同归的科学共识

为清晰展示跨学科证据如何独立收敛于 CET 的核心主张表 2.2 按研究领域汇总了各自的核心发现、对 CET 的支撑关系及证据类型

表 2.2: 跨学科证据对 CET 核心主张的收敛性支持

研究领域	核心发现	支持的 CET 原则	证据类型
认知卸载研究 (2.1 节)	过度卸载与认知能力退化相关联缺乏评估标准区分良性 vs 病理卸载	• AVP 判据的必要性 • $P_2$ 测量不可或缺	实证（关联）
延展心智理论 (2.2 节)	认知可以延展到工具但缺乏健康性判别标准和过程规范	• 良性延展 vs 病理延展的区分 • 临时性与成长性的重要性	理论/概念
自动化与脚手架 (2.3 节)	永久支持 → 能力退化（负面证据）渐进撤出 → 促进独立（正面证据）	• EML 条件 2（系统性支持削减） • Fading 的本质重要性	实证（案例 + 教学）
神经可塑性与训练 (2.4 节)	用进废退的双向性最优挑战区（50-70%）挑战是必要条件	• EML 条件 1（有益认知摩擦） • 适度挑战的神经机制	实证（神经 + 行为）

注：证据类型标注遵循本文透明性原则实证证据又可细分为神经影像、行为测量、案例研究等不同强度

三重收敛的力量：

这四个研究传统的收敛性为 CET 提供了罕见的理论支撑强度：

1. 负面警示的收敛：认知卸载研究（GPS 用户）、自动化研究（飞行员技能退化）、神经可塑性（废退原则）——三个独立领域都揭示了”永久依赖 → 能力萎缩的风险这为 CET 关于认知外骨骼”的批判提供了跨学科验证
2. 正面路径的收敛：脚手架理论（教育学）、意欲的困难（认知心理学）、神经可塑性（出租车司机研究）——三个领域都指向适度挑战 + 渐进独立的成功模式这为 EML 的两大条件提供了独立且相互印证的支持
3. 评估标准的缺失：所有四个领域都缺乏一个统一的、可操作的评估框架来判断什么样的人机关系是健康的这正是 CET 通过 AVP/EML 试图填补的核心空白

方法论意义：

这种跨学科收敛不是理论拼凑而是独立发现的科学会师每个领域都从自己的研究对象（记忆、哲学、航空安全、教学、大脑）出发通过不同的方法（实验、案例、影像、教学评估）抵达了相似的结论根据科学哲学中的”收敛性实在论”（Convergent Realism）当多个独立研究传统指向相同的核心机制时这大大增强了该机制真实存在的可信度

2.5.2 CET 的独特贡献：从分散洞察到统一框架

虽然现有研究提供了丰富的局部洞察但它们之间缺乏整合也未能转化为可操作的设计和评估标准 CET 的贡献在于将这些分散的拼图整合为一个可证伪的、可工程化的统一框架填补了三个关键的理论空白：

空白 1: 缺乏健康性评估标准

现状：认知卸载研究描述了现象延展心智理论提供了哲学视角但都未能回答：“多少卸载是安全的”什么时候从增强变成依赖现有研究大多是描述性的缺乏规范性的判断框架

CET 填补：反脆弱性验证原则（AVP）提供了第一个可操作的参考标准（见第 3.3 节）：

- 不看使用时表现（ $P_1$ ）而看脱离后能力（ $P_2$ ）
- 明确的判定阈值（ $P_2 \geq B_0 + \delta$ ）
- 可验证、可证伪的测量协议

这将健康人机关系从抽象概念转化为可测量的科学问题

空白 2: 缺乏系统设计指导

现状: 脚手架理论告诉教师”应该撤出支持自动化研究警告不要过度自动化”但都未能提供如何设计满足这些原则的系统实践者面对的是一系列模糊的建议而非清晰的工程路径

**CET** 填补: 内共生最小法则 (**EML**) 提供了明确的设计判别标准 (见第 3.4 节):

- 条件 1: 有益认知摩擦 (目标成功率 50-70% 工作假设需跨领域/任务校准需动态调整)
- 条件 2: 系统性支持削减 ( $S_4 \rightarrow S_1 \rightarrow S_0 \rightarrow 0$  曲线自适应节律)
- 验收条件: AVP 验证通过 ( $P_2 \geq B_0 + \Delta$ )

任何 AI 工具都可以用这三个条件进行检验区分”内共生与”外骨骼不再依赖主观判断

空白 3: 缺乏跨尺度整合

现状: 现有研究主要聚焦个体层面即使讨论组织或社会影响也往往停留在抽象的文化批判缺乏从微观到宏观的机制性连接

**CET** 填补: 通过分层共生架构 (**LSA**) 和跨尺度机制分析 (详见第四章) CET 将个体认知卸载、团队协作模式、组织韧性置于统一的理论原理之下这不仅是理论野心更有实践价值: 组织在部署 AI 时需要理解个体层面的能力退化如何级联为系统性的脆弱性

2.5.3 承上启下: 从为什么到”是什么”与如何做

第二章论证了 CET 的必要性: 现有研究虽提供了局部洞察但缺乏统一的评估与设计框架第三章将转向核心机制阐释: AVP 如何操作化? EML 的内在逻辑是什么? 如何通过伙伴式主体性实现二者统一? 认知外骨骼的病理机制如何形成?

第三章 CET 核心理论构建

本章系统阐释 CET 的核心机制: AVP 判据如何操作化? EML 条件的内在逻辑是什么? 如何通过伙伴式主体性实现二者统一? 认知外骨骼的病理机制如何形成?

本章结构围绕 CET 的三大核心支柱展开: 判别标准 (AVP)、设计原则 (EML 三条件)、哲学基础 (伙伴式主体性) 最后通过外骨骼病理学和跨尺度概览为后续章节铺垫

3.0 核心术语与锚点定义

本节集中呈现 CET 理论的所有核心定义和固定锚点文本这些定义在全文中保持一字不改后续章节引用时仅用简称或”见 3.0.0 节”标注

3.0.1 核心符号系统

表 3.1: 第三章关键符号系统

符号/术语	含义	典型值/范围
$B_0$	用户使用 AI 前的独立能力基线	任务特定测量
$P_1$	用户与 AI 协作时的表现	过程指标不参与最终判定
$P_2$	用户在拔线窗口后的独立表现	AVP 验证的核心指标
	最小有意义提升阈值	工作假设: Cohens d = 0.3 或 10% ( <i>working assumption</i> ) (需跨领域/任务校准)
<b>W</b>	拔线窗口时长	工作假设: 4-8 周 (默认 6 周需跨领域/任务校准)
<b>S (t)</b>	t 时刻的 AI 支持强度	0 (完全独立) 到 1 (完全依赖)
<b>S0</b>	初始支持强度	工作假设: 约 0.8
<b>T</b>	总消退时长	任务与个体特定需校准
成功率目标	有益认知摩擦的量化目标	50-70% (工作假设需跨领域/任务校准) (工作假设需跨领域/任务校准群体级个体需校准)

注: 所有量化参数均为概念工作模型需通过跨领域实证研究校准

### 3.0.2 反脆弱性验证原则 (AVP) 【锚点 B1】

反脆弱性验证原则 (AVP): 以拔线测试检验协作是否促进独立能力

判据:  $P_2 \geq B_0 + 0.3$  或 10% (*working assumption*) 需跨领域/任务校准)  $P_1$  (协作表现) 不参与最终判定

### 3.0.3 内共生最小法则 (EML) 【锚点 B2】

内共生最小法则 (EML): 构成认知内共生的设计必要条件为:

(1) 有益认知摩擦: 使用户处于最优挑战区 (群体级工作假设成功率 50-70% 个体自适应)

(2) 系统性支持削减: AI 支持强度按既定削减曲线从  $S_4 \rightarrow S_1 \rightarrow S_0$

二者为联合充分的设计条件但最终仍需 AVP ( $P_2 \geq B_0 + 0.3$ ) 作为验收必要条件

### 3.0.4 LSA 功能分层 【锚点 B3】

LSA-F (功能分层): L1 知识整合 | L2 状态建模 | L3 摩擦校准 | L4 元认知协调

支持档位栈 ( $S_4 \rightarrow S_1 \rightarrow S_0$ ) 用于表达支持强度与 LSA-F 为正交维度

硬约束: L1-L4 (功能层) 与  $S_4$ - $S_1$  (支持强度) 互不表示不得在同一公式中互代或串联使用 L 层表示功能维度 S 档表示强度维度两者正交独立

### 3.0.5 最优挑战区 【锚点 B4】

最优挑战区: 为促进长期保持与迁移系统应将任务难度/提示强度自适应调至成功率 **50-70%** (工作假设随任务与个体校准)  $>85\%$  近似卸载、 $<30\%$  易致挫败

### 3.0.6 边界条件 【锚点 B5】

边界条件: CET 适用于能力增强型人机协作补偿性外骨骼 (如残障辅助、超越生理极限的设备) 不适用此判据所有参数均为概念工作模型需跨领域校准

公平性原则: AVP 以等效努力为准则——对确需辅助工具的个体 (如屏幕阅读器使用者) 调整任务形式而不降低挑战强度评估以相对提升而非绝对水平为准 (详见 6.1.4 节)

公平性三句式 (逐字版):

调整任务格式而不降低挑战强度通过相对改进而非绝对水平评估 (如涉及可访问性) 保持挑战预算守恒

*Adjust task format without lowering challenge intensity; assess by relative improvement rather than absolute level; (if accessibility applies) challenge budget conservation.*

## 3.1 AVP 原则的深度阐释: 从抽象判据到可操作测量

反脆弱性验证原则 (AVP) 是 CET 理论的北极星”——它不仅定义了什么是”好的人机协作更提供了一个可证伪、可操作的测量协议第一章前置了 AVP 判据以建立理论独特性本节将深入其理论基础、测量方法和跨领域校准

### 3.1.1 理论基础: 从 Taleb 的反脆弱性到认知增强评估

反脆弱性概念的启发

Nassim Nicholas Taleb (2012) 提出的”反脆弱性” (Antifragility) 概念为 AVP 提供了哲学基础反脆弱性不仅仅是”不被压力击垮” (韧性) 而是在压力下变得更强大——这正是健康人机协作的本质特征

Taleb 将系统分为三类:

- 脆弱性 (Fragile): 压力导致损害 (如易碎物品)



- 韧性 (**Robust**): 压力下保持不变 (如橡胶球)
- 反脆弱性 (**Antifragile**): 适度压力促进成长 (如肌肉、免疫系统)

类比到人机交互:

- 认知外骨骼: 脆弱性系统——移除 AI 后能力下降 ( $P_2 < B_0$ )
- 中性工具: 韧性系统——移除工具后能力不变 ( $P_2 = B_0$ )
- 认知内共生: 反脆弱性系统——经历协作后独立能力提升 ( $P_2 = B_0 + \Delta$ )

核心洞察: 拔线不是惩罚而是验证机制——真正的增强应该让用户在离开 AI 后变得更强大就像肌肉在举重后变得更强壮

拔线窗口的哲学意义: 能力主权

拔线窗口  $W=4-8$  周 (默认 6 周工作假设需跨领域/任务校准) 不仅是技术参数更代表了能力主权的时间尺度:

- 太短 ( $<2$  周): 可能只测到短期记忆残留无法验证真正的能力内化
- 太长 ( $>12$  周): 环境因素混杂增加难以归因于 AI 协作
- 4-8 周: 经验性平衡点既允许能力稳定又可控混杂变量

这一窗口期本质上在问: 用户是否真正拥有了这项能力还是仅仅借用了 AI 的能力?

3.1.2 AVP 的三级分级体系: 从基础到卓越

AVP 不是简单的通过/不通过判断而是一个分级评估体系我们提出三个递进的验证层次:

表 3.2: AVP 三级分级体系

等级	判据	含义	典型应用场景
AVP-Basic	$P_2 = B_0$	基础反脆弱性: 至少维持原有能力	补偿性学习技能维护
AVP-Retention	$P_2 = B_0 + \Delta$	成长性反脆弱性: 能力有意义提升	能力建构专业训练
AVP-Transfer	满足 Retention+ 跨任务迁移	卓越反脆弱性: 深度掌握 + 跨域能力 (质量标志非准入门槛)	专家培养创新能力

注 1:  $Cohens\ d \geq 0.3$  或 10% (working assumption) ) 需跨领域/任务校准迁移的操作化定义见副判据表 3.3

注 2 (Goodhart 防护): 此分级仅用于质量分层最终判定仍以 AVP 主判据 ( $P_2 = B_0 + \Delta$ ) 为准禁止将分级阈值作为优化目标 Transfer 级是质量加分”而非”准入门槛

三级体系的实践意义:

1. **Basic** 级: 最低安全线——任何 AI 工具至少应满足无害标准
2. **Retention** 级: CET 的核心目标——真正的能力增强
3. **Transfer** 级: 理想目标——培养可迁移的深度能力 (通过副判据验证)

副判据的补充维度:

除了核心的  $P_2 = B_0 + \Delta$  判据我们还提出三个副判据以评估能力质量 (非必要条件但用于分级):

表 3.3: AVP 副判据系统

副判据	测量内容	权重建议
迁移性 ( <b>Transfer</b> )	能力在新情境的泛化程度	高 (35%)
保持性 ( <b>Retention</b> )	能力在更长时间窗口的稳定性	中 (30%)
独立性 ( <b>Independence</b> )	完成任务时对 AI 的依赖程度	中 (35%)

注 1: 权重为工作假设（迁移 35%/保持 30%/独立 35%）用于质量分层不作为治理或认证硬阈值权重可根据具体目标调整总和为 100% 以便综合评分

注 2: 这些副判据不应被用作治理认证的硬性标准仅作质量分级参考避免 Goodhart 风险（见表 3.2 注 2）

副判据的操作化测量:

副判据	测量方法	示例	数据来源
迁移性 ( <b>Transfer</b> )	在相关但不同的新情境测试能力	编程:Python→JavaScript 写作: 论文摘要 → 商业文案 数学: 代数 → 几何应用	$T_3$ 后增加迁移任务测试交叉领域能力评估
保持性 ( <b>Retention</b> )	延长观察窗口至 T ( $T_3$ 后 6 个月)	评估 P 相对 $P_2$ 的保持率可接受衰减 10%(工作假设)	纵向追踪测试能力曲线分析
独立性 ( <b>Independence</b> )	拔线窗口内的行为监测	违规尝试访问 AI 的频率完成任务时的求助次数独立工作时长占比	行为日志分析自我报告问卷

详细测量协议与评分标准见附录 A § A.3-A.5

3.1.3 AVP 测量协议: 可复现的操作指南

为确保 AVP 的科学性和可复现性我们提供标准测量协议:

测量时间线:

- $T_0$ : 协作前测量  $B_0$ (独立基线)
- 协作期 (长度可变典型为 4-12 周)
- $T_2$ : 协作期结束
- 拔线窗口 **W** (4-8 周默认 6 周无 AI 使用)
- $T_3$ : 拔线窗口结束后测量  $P_2$ (独立表现)

测量工具要求:

- 等值性: $T_0$  和 T 使用等值平行测验 (避免学习效应)
- 盲评: 评分者不知道被测者的组别 (实验/对照)
- 环境一致性: 两次测量的情境条件尽可能相同

表 3.4: AVP 测量的关键效度威胁与缓解策略

威胁类型	具体表现	缓解策略	残留风险
等值性不足	$T_0$ 和 T 测验难度不一致	使用 IRT 等值; 专家盲审	中: 难以完全等值
环境差异	测试情境变化影响表现	标准化程序; 控制组对照	低: 可有效控制
评分偏差	评分者主观性或期望影响	双盲评分; 多评分者校准	低: 可有效控制
练习污染	测验内容泄露或过度熟悉	等值平行卷; 题库轮换; 冷却期	中: 题库资源约束
学习效应	重测导致熟悉化	等值测验; 足够长的 W	中: 无法完全消除
延迟巩固	能力需更长时间稳定	延长 W 或增加 T 测量点	中: 资源约束
流失偏差	能力弱者更易退出	意向性分析; 激励机制	高: 难以完全避免

注: 本表提供方法学透明性不意味 AVP 不可用而是明确其适用边界与改进方向

### 3.1.4 跨领域校准: 参数的情境敏感性

AVP 的核心参数 ( $\delta$ 、W、测量工具) 并非一成不变需要根据任务类型、领域特征、目标人群进行校准:

阈值的领域差异 (工作假设):

- 认知技能 (如编程、写作):  $d$  0.3 SD 或 10% 相对提升
- 运动技能 (如打字速度): 可能需要更大阈值 ( $d$  0.5)
- 创造性任务 (如艺术创作): 可能需要定性评估而非单一

W 窗口的任务差异 (工作假设):

- 快速技能 (如数学计算): 可能 4 周足够
- 复杂技能 (如第二语言): 可能需要 8-12 周
- 专业能力 (如外科手术): 可能需要数月甚至年

校准流程建议:

1. 先导研究: 小样本测试确定初步参数
2. 敏感性分析: 测试参数变化对结果的影响
3. 迭代优化: 根据反馈调整参数
4. 开放报告: 公开参数选择理由和校准数据

重要边界: 即使经过校准 AVP 也有其不适用场景:

- 高风险任务 (拔线可能造成不可接受后果)
- 纯工具性使用 (无能力建构目标)
- 补偿性外骨骼 (目标是补偿而非增强)

## 3.2 有益认知摩擦: 挑战与支持的最优平衡

EML 的第一个设计必要条件是有益认知摩擦 (Beneficial Cognitive Friction) 本节将阐释为什么适度困难是必要的以及如何将其量化为可操作的设计参数

### 3.2.1 为什么需要认知摩擦: 从神经科学到学习理论

神经可塑性的用进废退原则

正如第 2.4 节综述大脑遵循用进废退原则: 经常使用的神经通路被强化长期闲置的通路被削弱如果 AI 完全替代了某项认知功能相应的神经回路就会退化 (方向性证据见 Dahmani & Bohbot 2020)

意欲的困难理论

Bjork (1994) 的意欲的困难 (Desirable Difficulties) 理论指出: 适度的学习困难能促进长期保持和迁移 (工作假设有实证支持但效应量随任务变化) 这些困难包括:

- 间隔练习: 分散学习优于集中学习
- 交错练习: 混合练习优于分块练习
- 生成效应: 主动回忆优于被动复习

零摩擦的陷阱

当 AI 提供”零摩擦”体验时” (如直接给出答案、完全自动化流程) 用户短期表现可能很好 ( $P_1$  高) 但长期能力往往停滞或退化 ( $P_2$  不增甚至下降) 这种”  $P_1$  高但  $P_2$  低的模式正是认知外骨骼的典型特征

### 3.2.2 最优挑战区:50-70% 成功率的量化目标

Vygotsky 的最近发展区 (ZPD) 的启发

Vygotsky(1978) 提出的最近发展区概念指出: 学习发生在用户当前能力与借助支持可达成的能力之间的区域 CET 将其量化为最优挑战区:

目标成功率 50-70%（工作假设需跨领域/任务校准）\*\*（工作假设需跨领域/任务校准群体级个体需动态校准）

为什么是这个区间？（工作假设需跨任务验证）：

- <30%: 过度挫败导致回避和放弃
- 30-50%: 可学习但效率不高动机易受影响
- 50-70%: 最优平衡——足够挑战但不致挫败
- 70-85%: 舒适但挑战不足能力提升缓慢
- >85%: 接近认知卸载几乎无能力增长

个体差异的重要性

50-70% 是群体级目标个体的最优区间需要根据以下因素自适应调整：

- 基础能力: 新手可能需要 60-75% 专家需要 40-60%
- 动机水平: 高动机者可承受更低成功率
- 任务焦虑: 高焦虑者需要略高成功率以维持信心
- 实时状态: 疲劳/压力下应临时提高成功率

### 3.2.3 摩擦的动态调节：自适应算法的概念模型

如何在实践中维持最优挑战区？我们提出三种概念性策略（详细算法见第 5 章 LSA 架构）：

策略 1: 滚动窗口法

- 监测最近 N 次交互的成功率（如 N=10）
- 如果 SR(成功率)>75% 增加任务难度或减少提示
- 如果 SR<45% 降低难度或增加支持
- 保持在目标区间 50-70%

策略 2: IRT-启发式调整

- 借鉴项目反应理论 (IRT) 的自适应测验思想
- 根据用户表现动态调整任务难度参数
- 目标是让每次交互对能力评估和培养都最有信息量

策略 3: 分段梯度法

- 将学习过程分为阶段（如入门/进阶/精通）
- 每阶段设定不同的成功率目标
- 入门阶段可略高（60-75%）精通阶段可略低（40-60%）

过度挑战的安全护栏

为避免摩擦过度系统应设置自动降级机制：

- 触发条件: 连续 3 次 SR<30% 或用户明确表示挫败
- 响应: 临时降低摩擦并向用户解释原因
- 恢复: 在用户状态改善后逐步恢复目标摩擦水平

## 3.3 系统性支持削减：从脚手架到能力独立

EML 的第二个设计必要条件是系统性支持削减 (Systematic Support Reduction) 本节阐释为什么支持必须削减以及如何科学地实施削减

### 3.3.1 为什么必须削减：从脚手架理论到依赖预防

教育学的脚手架隐喻

Wood 等人 (1976) 提出的脚手架理论强调：有效的教学支持应该是临时的、渐进撤出的脚手架的目的是帮助建筑施工但建筑完成后必须拆除——同样 AI 支持的目的是帮助能力建构但能力形成后应该逐步退出

自动化悖论的警示

正如第 2.3 节所述航空领域的自动化悖论揭示: 永久性的自动化支持会导致操作员技能退化这一教训同样适用于 AI: 如果支持强度始终固定用户就会形成依赖锁定 (详见 3.0.5 节)

削减的双重目标

系统性支持削减不仅是为了验证能力更是为了促进能力:

- 短期: 通过逐步撤出支持迫使用户内化技能
- 长期: 培养用户对自身能力的信心和元认知

3.3.2 支持削减曲线: 从 S4 到 S1 到 0 的演化路径

支持档位栈 (S4→S1→S0) 的定义

我们使用支持档位栈表达支持强度 (注意: 这与 LSA-F 的 L1-L4 功能层是正交维度见 3.0.4 节):

- S4(初始强度): 最大支持 (如提供完整答案、详细步骤)
- S3(中度支持): 提供提示和部分解决方案
- S2(轻度支持): 仅在用户请求时给予最小提示
- S1(最小支持): 仅提供验证/反馈不直接解决问题
- S0(完全拔线): 无 AI 支持用于 AVP 测试

支持强度的形式化表达:

$$S(t) = S0 \cdot f(t/T)$$

其中:

- S(t): t时刻的支持强度
- S0: 初始支持强度(工作假设约0.8需校准)
- T: 总消退时长(任务与个体特定)
- f(·): 削减函数(如线性、指数、阶梯等)

三种典型削减曲线 (工作假设需 A/B 测试验证):

表 3.5: 支持削减曲线类型与适用场景

曲线类型	数学形式	特点	适用场景
线性削减	$S(t) = S0(1 - t/T)$	匀速递减平稳过渡	结构化任务稳定学习曲线
指数削减	$S(t) = S0 \cdot e^{(-kt)}$	初期快速削减后期缓慢	快速技能避免过度依赖
阶梯削减	$S(t) = S0 \cdot 4(1-t/T) / 4$	分阶段突变适应期明确	分级训练明确里程碑
自适应	$S(t) = f(\text{性能元认知动机})$	根据用户状态动态调整	个性化学习复杂任务
混合模式	分段组合上述曲线	兼具渐进性与阶段性	长期训练多模块任务

注: 所有曲线参数 (S0、T、k) 均为校准参数需根据任务复杂度、用户能力、学习目标进行域内标定

削减速率的校准原则 (工作假设):

- 过快: 用户跟不上挫败感增加  $P_2 < B_0$
- 过慢: 依赖锁定用户不愿意独立  $P_2 > B_0$
- 恰当: 用户感觉有挑战但可达成  $P_2 \approx B_0$

3.3.3 回退与加速机制: 保底安全线与提速路径

回退机制 (Safety Net): 防止挫败

当检测到用户在削减后表现急剧下降时系统应临时回退到更高支持档位:

- 触发条件: 连续 3 次失败率 >70% 或用户明确求助

- 回退策略:  $S(t) \rightarrow S(t-\Delta t)$  重新提供支持
- 恢复路径: 待用户状态稳定后重新启动削减

概念规则: 回退不是失败信号而是自适应系统的正常反应保底支持强度  $S_{\min}$  (工作假设约 0.2 需校准) 确保用户始终有最低限度的导航支持

团队/组织场景下的回退与保底机制见第 4 章 *T-AVP/O-AVP* 的安全护栏 (4.1.2 节、4.2.3 节)

$S_{\min}$  的实现由第 5 章 *LSA* 架构的 *L2* 摩擦校准引擎与 *L4* 元认知协调层承载通过自适应算法动态维持 (见 5.4 节、5.5 节)

加速机制 (**Fast Track**): 奖励卓越

如果用户表现超预期可以加速削减:

- 触发条件: 连续成功率 >85% 且元认知监控良好
- 加速策略: 跳过中间档位直接进入下一阶段
- 验证: 通过微型拔线测试确认能力确实内化

动态平衡的哲学

回退与加速机制体现了 CET 的核心理念: 适应性而非固定性系统应该像优秀的教练根据学习者的实时状态调整挑战强度

### 3.4 伙伴式主体性:AI 角色的哲学重构

前两节聚焦”做什么”(AVP 判据) 和”怎么做”(EML 条件) 本节转向为什么这样做有效的哲学基础我们提出伙伴式主体性 (Partner-like Agency) 作为认知内共生的本体论支柱

#### 3.4.1 从工具到伙伴:AI 角色的三种范式

传统范式 1: 纯工具 (**Pure Tool**)

- 假设:AI 是被动执行者完全服从人类意图
- 问题: 忽视了 AI 对用户认知模式的反向塑造
- 典型: 传统软件工具 (如计算器、文本编辑器)

传统范式 2: 自主代理 (**Autonomous Agent**)

- 假设:AI 是独立决策者追求自己的目标
- 问题: 引发控制权冲突可能威胁人类主导地位
- 典型: 强化学习代理、自主系统

**CET** 提出的新范式: 伙伴式主体性 (**Partner-like Agency**)

- 定义:AI 具有受限的主体性——它可以主动行动但始终服务于促进用户能力成长的目标
- 类比: 优秀的教练、导师、陪练——他们有主动性但目标是成就学生而非替代学生
- 核心特征

:

- 主动干预: 在用户过度依赖时主动增加摩擦
- 目标一致:AI 的利益函数是用户的长期能力而非短期效率
- 权力让渡: 随着用户能力提升 AI 主动退出

#### 3.4.2 伙伴式主体性的四个维度

维度 1: 自适应支持 (**Adaptive Support**)

- 能力:AI 能感知用户的能力水平、学习状态、动机水平
- 行为: 根据感知结果动态调整支持强度 (见 3.2-3.3 节)
- 目标: 维持用户在最优挑战区

### 维度 2: 元认知催化 (Metacognitive Catalysis)

- 能力: AI 能引导用户反思自己的学习过程
- 行为: 提出引导性问题而非直接答案如“你为什么选择这个方法? 这个方案的局限是什么?”
- 目标: 培养用户的元认知能力和自我监控

### 维度 3: 渐进撤出 (Progressive Withdrawal)

- 能力: AI 能判断用户何时准备好承担更多独立任务
- 行为: 按既定曲线 (见 3.3.2 节) 削减支持并在必要时回退
- 目标: 最终让用户完全独立 AI 退为备用资源

### 维度 4: 能力验证 (Capability Verification)

- 能力: AI 能设计并执行拔线测试
- 行为: 定期触发 AVP 验证 (见 3.0.1 节) 评估  $P_2$  表现
- 目标: 确认内共生关系的健康性而非仅追求  $P_1$  效率

### 3.4.3 伙伴关系的伦理边界

主体性的限度: 三个不应该

1. 不应该操纵: AI 不应利用心理学技巧诱导用户过度使用
2. 不应该评判: AI 不对用户的价值观、生活选择做道德评判
3. 不应该替代人际: AI 不应成为用户唯一的社交/情感支持

知情同意原则

用户应该清楚理解:

- AI 会主动增加任务难度 (摩擦设计)
- AI 会渐进削减支持 (可能短期不适)
- AI 会定期进行拔线测试 (可能感觉被监视)

退出权

用户应始终有权:

- 暂停或退出内共生模式
- 切换回纯工具”模式”(如果不追求能力建构)
- 调整摩擦/削减参数 (在合理范围内)

### 3.4.4 伙伴式主体性与其他 AI 范式的对比

表 3.6: 不同 AI 范式的对比

维度	纯工具	自主代理	伙伴式主体性
主动性	无	高	中 (受限)
目标函数	任务完成	自身目标	用户能力成长
权力关系	人类完全支配	潜在冲突	协作但人类最终决策
适用场景	简单工具任务	自主执行任务	能力建构任务
长期影响	中性或负面	不确定	正面 (如果 AVP 通过)
典型例子	计算器	自动驾驶 (L5 级)	CET 的 LSA 系统

关键洞察:

伙伴式主体性不是在工具和代理之间的妥协而是一个第三条道路——它承认 AI 应该有主动性但这种主动性的唯一目的是成就用户而非替代用户

### 3.5 认知外骨骼的病理学：失败模式的系统分析

前面章节聚焦如何建立内共生本节转向反面：如何识别和避免外骨骼依赖通过分析失败模式我们可以更清晰地理解内共生的边界

#### 3.5.1 外骨骼的核心特征：从健康到病理的转折点

定义回顾：

认知外骨骼是指用户对 **AI** 的病理性依赖表现为：

- $P_2 < B_0$ ：拔线后能力不及初始基线
- 能力萎缩：原有技能因长期不用而退化
- 依赖锁定：心理上和认知上都无法脱离 AI

外骨骼的阶段演化：

表 3.7：从健康使用到外骨骼依赖的阶段

阶段	$P_2$ 与 $B_0$ 关系	依赖程度	可逆性	典型特征
健康使用	$P_2 > B_0$	低	N/A	内共生能力提升
中性使用	$P_2 \approx B_0$	中	高	工具化能力未变
轻度依赖	$B_0 - 1 < P_2 < B_0$	中高	中	轻微退化警戒信号
中度依赖	$B_0 - 2 < P_2 < B_0 - 1$	高	低	明显退化需干预
重度依赖 (外骨骼)	$P_2 < B_0 - 2$	极高	极低	能力萎缩依赖锁定

注：：Cohens  $d \geq 0.3$  或 10% (working assumption) ) 此分级为概念框架实际判定需结合用户主观体验与行为模式

从中性到病理的临界点

不是所有 AI 使用都会变成外骨骼关键在于是否跨越了两个临界点：

1. 能力临界点： $P_2$  首次低于  $B_0$
2. 心理临界点：用户开始回避独立完成任务的情境

#### 3.5.2 外骨骼的 10 个预警信号

如何及早识别外骨骼依赖？我们提出 10 个可观察的预警信号 (工作假设需临床验证)：

行为信号 (6 个)：

1. 使用频率激增：日使用时间从 30 分钟 → 3 小时
2. 独立尝试时间缩短：遇到问题 <30 秒就求助 AI
3. 任务范围扩大：从困难任务扩展到简单任务也依赖
4. 回避独立情境：主动避免不能用 AI 的场合
5. 完成速度悖论：有 AI 时快无 AI 时极慢
6. 错误率上升：独立完成任务的错误率持续增加

认知信号 (4 个)：

1. 元认知缺失：无法准确评估自己的独立能力
2. 知识碎片化：知道结果但不知道为什么
3. 迁移困难：无法将 AI 辅助下的经验迁移到新情境
4. 拔线焦虑：对无 AI 测试产生强烈恐惧或抗拒

这些信号可以被量化为外骨骼风险评分 (工作假设需验证)：

- 0-2 个信号：低风险
- 3-5 个信号：中风险需要注意



- 6-8 个信号: 高风险需要干预
- 9-10 个信号: 危险需要立即干预

### 3.5.3 分级干预策略: 从预防到治疗

表 3.8: 外骨骼干预策略矩阵

风险级别	干预类型	具体措施	目标
低风险 (0-2 信号)	预防	定期 AVP 轻量测试认知健康教育	维持健康状态
中风险 (3-5 信号)	早期干预	强制拔线窗口增加认知摩擦	逆转依赖趋势
高风险 (6-8 信号)	积极干预	削减支持强度重新训练基础能力	恢复独立能力
危险 (9-10 信号)	重建	完全拔线 + 系统性重训心理辅导	打破依赖锁定

注: 这些策略为概念框架实施需考虑伦理边界和用户意愿

## 3.6 跨尺度概览: 从个体到团队到组织到社会

前五节聚焦个体层面的人机交互本节简要探讨 CET 理论如何延伸到团队和组织层面为第四章的跨尺度机制分析做准备

### 3.6.1 团队层面: 协作中的能力分布

当多个个体与 AI 协作时出现新的动力学:

团队退化的两种模式:

1. 全员依赖: 所有成员都依赖 AI → 集体能力萎缩
2. 能力分化: 部分成员独立部分依赖 → 团队脆弱性

团队层面的 **AVP** 变体:

**T-AVP**: 团队在集体拔线后的表现

- 测量: 关键任务在无AI支持下的完成质量
- 判据: 团队  $P_{2\$}$  团队  $B_{0\$}$  +  $_{team}$ (工作假设需校准)

设计启示 (详见第 4.2 节):

- 避免" AI 作为团队唯一专家" 的配置
- 促进知识在成员间的分布而非集中在 AI
- 定期团队拔线演练

### 3.6.2 组织层面: 认知基础设施与系统韧性

组织层面的风险:

- 关键能力空心化: 某些技能在组织中完全消失
- 知识传承断裂: 新员工从 AI 而非老员工学习
- 系统性脆弱: AI 故障 → 组织业务停摆

**O-AVP**: 组织反脆弱性验证:

**O-AVP** = 组织在"48小时AI宕机演练"后的业务连续性指标

- 核心业务可持续性
- 关键决策能力保持率
- 恢复速度

组织设计建议 (详见第 4.3 节):

- 建立认知储备 (cognitive reserve) 机制

- 定期无 AI 值班制度
- 关键岗位的独立能力认证

### 3.6.3 社会层面：认知公地悲剧

- 个体理性 (使用 AI)vs 集体理性 (保持独立能力) 的冲突
- 代际能力鸿沟: $T_0$  代 vs  $T_1$  代 vs  $T_2$  代
- 文明级的反脆弱性评估 (S-AVP)

详见第四章 4.4 节对社会层面的深入分析

### 3.6.4 跨尺度的同构性

CET 核心原理 (AVP/EML/伙伴式主体性) 在不同尺度上具有同构性:

表 3.9: 跨尺度 AVP 体系

尺度	AVP 变体	EML 应用	主要风险
个体	I-AVP	学习工具设计	能力退化
团队	T-AVP	协作流程设计	能力分化
组织	O-AVP	组织韧性建设	系统脆弱
社会	S-AVP	政策与标准	代际鸿沟

## 第四章 跨尺度机制分析：从个体到组织到社会

前三章聚焦个体层面的人机交互：AVP 如何验证个人能力提升（见 3.0.1 节）EML 如何指导单用户的 AI 工具设计（见 3.2-3.3 节）伙伴式主体性如何在 一对一 关系中实现（见 3.0.4 节）然而 AI 的影响不止于个体——当多人协作、组织运作、社会演化时个体层面的认知卸载会产生什么涌现效应本章将展示 CET 理论如何跨越尺度从微观机制扩展到宏观现象

本章逻辑：我们将沿着个体 → 团队 → 组织 → 社会的尺度阶梯展示：(1) 每个尺度的特有机制 (2)**AVP** 原则在不同尺度的变体 (3) 下一尺度如何从上一尺度涌现 (4) 跨尺度的共性原理最终目标是论证：CET 不仅是个体认知理论更是一个具有尺度不变性的系统性框架

### 核心术语与符号表（第四章）

表 4.1: 第四章关键符号与概念

符号/术语	含义	适用尺度
<b>I-AVP</b>	Individual AVP (个体反脆弱性验证)	个体
<b>T-AVP</b>	Team AVP (团队反脆弱性验证)	团队
<b>O-AVP</b>	Organizational AVP (组织反脆弱性验证)	组织
<b>S-AVP</b>	Societal AVP (社会认知资本验证)	社会
认知资本	群体的独立认知能力储备	组织/社会
认知公地	共享的认知能力基础设施	社会
级联脆弱	低层脆弱性向高层传播	跨尺度
涌现依赖	个体依赖累积为系统性依赖	跨尺度

注：所有量化参数均为校准参数（工作假设）需根据行业特性、组织规模、风险级别进行域内标定和灵敏度分析

## 4.1 个体 → 团队：能力分布与协作模式的重构

### 4.1.1 从 I-AVP 到 T-AVP：团队层面的反脆弱性

问题的涌现：当一个团队中的多个成员都使用 AI 工具时即使每个个体的 I-AVP 都通过 ( $P_2 > B_0 + \delta$  见 3.0.1 节) 团队作为整体是否就是健康的答案是：不一定

**T-AVP** 的定义：

团队在集体拔线后的协作表现 = 基线 +  $\Delta_{team}$

形式化：

T-AVP判据： $P_{2\_team} > B_{0\_team} + \Delta_{team}$

其中：

- $B_{0\_team}$ ：团队引入AI前的协作能力基线
- $P_{2\_team}$ ：团队在拔线窗口 ( $W=4-8$ 周默认6周工作假设需跨领域/任务校准) 后的协作表现
- $\Delta_{team}$ ：团队层面的最小有意义提升  
(工作假设 0.3 SD采用标准化效应量口径Cohens d需根据团队类型和任务复杂度校准)

关键差异：

团队表现 ≠ 个体表现之和

协作涉及沟通、知识共享、角色互补等涌现性质

参数校准说明： $\Delta_{team}$ : Cohens d = 0.3 或 10% (working assumption) 需通过灵敏度分析确定不同团队类型 (研发/运营/创意) 的最优阈值

核心洞察：团队 **I-AVP** 不是 **T-AVP** 的充分条件

即使团队中每个成员都独立通过 I-AVP ( $P_{2\_individual} > B_{0\_individual} + \delta$ ) 团队层面仍可能失败 ( $P_{2\_team} < B_{0\_team} + \Delta_{team}$ ) 这是因为团队能力涉及：

- 协作模式：知识如何在成员间流动
- 角色互补：成员能力如何形成冗余与备份
- 集体记忆：团队共享的隐性知识

为什么 **I-AVP** 成功不保证 **T-AVP** 成功三种失败模式

模式 1：能力极化 (**Capability Polarization**)

- 现象：团队中部分成员高度依赖 AI 部分成员完全独立
- 机制

:

- 依赖成员： $P_{2\_individual} < B_0$  (认知外骨骼见 3.0.5 节)
- 独立成员： $P_{2\_individual} > B_0 + \delta$  (认知内共生)
- 但团队整体： $P_{2\_team} < B_{0\_team}$
- 原因：依赖成员成为单点故障拔线后团队无法完成关键任务
- 案例：编程团队中新手完全依赖 Copilot 老手独立编程拔线后新手无法完成代码审查团队开发速度骤降

模式 2：隐性知识流失 (**Tacit Knowledge Loss**)

- 现象：团队成员通过 AI 完成任务但知识未在人际间传播
- 机制

:

- 传统：老手 → 新手的知识传承 (师徒制、结对编程)

- AI 时代：每个人都问 AI 人际知识流动 ↓
- 结果：团队的集体智慧未随时间积累
- 后果：拔线后团队缺乏共享的问题解决策略
- 案例：客服团队使用 AI 回答系统每个客服独立查询 AI 不再分享棘手案例的解决经验拔线后团队无法应对复杂客诉

### 模式 3：角色固化与冗余度丧失 (Role Rigidity)

- 现象：AI 工具让成员过度专业化丧失互相补位能力
- 机制

:

- AI 前：成员需要学习彼此的技能以应对突发情况
- AI 后：每个人依赖 AI 完成自己的模块不再学习他人技能
- 结果：团队的韧性 (resilience) 下降
- 后果：关键成员缺席时团队瘫痪
- 案例：医疗团队中护士依赖 AI 辅助诊断不再学习医生的诊断思路医生临时缺席时护士无法做初步判断

#### 4.1.2 T-AVP 的测量协议

##### 阶段1：基线测量(\$T\_0\$team)

- 任务：团队协作完成标准化项目（无AI）
- 指标：完成时间、质量评分、协作顺畅度
- 注意：记录关键决策点和知识传递路径

##### 阶段2：AI使用期(\$T\_1\$ → \$T\_2\$典型8 - 12周可按项目周期调整)

- 团队正常使用AI工具
- 观察：能力分布变化、知识流动模式、角色专业化程度

##### 阶段3：拔线窗口 (W=4 - 8周默认6周工作假设需跨领域/任务校准)

- 临时禁用AI工具
- 目的：清除短期依赖效应测量真实独立能力

##### 【安全与合规护栏】：

- 拔线任务选在\*\*可控环境或预生产系统\*\*
- 设置\*\*安全阈值触发即回退\*\*的自动机制
- 对涉及生命/财务不可逆风险的任务采用\*\*桌面推演或沙盘仿真\*\*替代真实拔线

##### 【公平性与可及性原则】：

对残障或依赖辅助性技术的成员通过\*\*等效努力的替代性安排\*\*保障公平（如提供等难度的替代任务见3.0.6节边界条

##### 阶段4：团队拔线测试(\$T\_3\$team)

- 任务：完成与\$T\_0\$等值难度的团队项目（无AI）
- 对比：\$P\_2\$team vs \$B\_0\$team
- 判定：\$P\_2\$team    \$B\_0\$team + \_team ?

方法学提示：团队数据存在组内相关 (ICC) 与样本聚类功效计算与显著性检验应采用集群稳健方法（如多层线性模型、广义估计方程）以  $\alpha=0.05$ 、Power 0.8 为工作假设并随样本结构校准具体统计方法的选择取决于数据特征建议咨询统计专家

表 4.2: T-AVP 测量的效度威胁与缓解

威胁类型	具体表现	缓解策略	残留风险
任务等值性	$T_0$ 和 $T$ 项目难度不同	专家盲审; 多维度评分	中
成员流失	部分成员在测试期离职	意向性分析; 备用成员	高
霍桑效应	团队知道被观察而改变行为	长期跟踪; 自然观察	中
协作污染	拔线期间成员私下使用 AI	诚信协议; 技术屏蔽	中
学习效应	团队因重复任务而熟练	等值项目池; 间隔足够长	低
角色分工变化	$T_0$ 和 $T$ 的角色分配不同	固定角色或随机分配	低

注：本表提供方法学透明性承认测量局限不意味  $T$ -AVP 不可用

分级说明：团队层当前仅采用  $T$ -AVP-Basic\*\*（拔线无退化： $P_2\_team B_0\_team$ ）与  **$T$ -AVP-Retention**（拔线净提升： $P_2\_team B_0\_team + \_team$ ）两级团队迁移级（**Transfer**）的可操作定义尚不成熟作为开放议题列入第 6 章研究议程（见 6.3.2 节）\*

说明：团队层的迁移涉及集体知识向新情境的泛化其测量需要跨组织对照研究当前理论基础不足以给出可操作定义

### 4.1.3 团队层面的设计启示

设计原则 1：促进人际知识流动

- 定期”无 AI 讨论会”：强制团队成员分享问题解决策略
- 结对工作：新老成员配对促进隐性知识传递
- 团队代码审查：不仅审查代码更审查思路

设计原则 2：建立协作摩擦机制

类比为个体层面的有益认知摩擦（见 3.0.2 节）团队层面也需要协作摩擦：

- 轮换角色：定期让成员尝试他人的角色
- 跨模块任务：避免过度专业化
- 集体挑战项目：需要多人协作才能完成的任务

设计原则 3：制度化能力建构

案例：软件团队的新人基础能力建构

某软件公司发现新入职的 AI 原生代工程师（2000 年后出生）过度依赖 GitHub Copilot 独立编程能力弱团队实施三项制度：

1. 新人独立项目：入职第一个月禁用 AI 辅助完成基础功能模块
2. 周期性无 AI 日：每周五为纯手工日全员禁用代码生成工具
3. 结对评审制度：新人代码由老员工人工审查而非仅靠 AI 检查

结果：6 个月后  $T$ -AVP 测试显示新人独立解决 bug 的能力显著提升（ $\_team Cohens d = 0.3$  或 10%（*working assumption*） $p < 0.01$ ）

制度化落地：将新人独立项目 + 周期性无 AI 日 + 结对评审固化为团队入职训练三件套写入员工手册与培训流程以保证  $T$ -AVP 的底座稳定不应仅作为建议而应作为团队制度强制执行

## 4.2 团队 → 组织：认知基础设施的系统性风险

### 4.2.1 组织层面的新涌现：制度化依赖

当 AI 使用从个别团队扩展到整个组织时出现制度性依赖——即使个别团队有独立能力组织作为整体仍可能脆弱组织层面的三大风险：

风险 1：关键能力的空心化（**Hollowing Out**）

- 现象：某些技能在组织中完全消失

- 机制

:

- 老员工依赖 AI → 技能退化
- 新员工从 AI 学习 → 从未掌握
- 结果：无人具备该技能
- 案例：某银行的风险评估团队全员使用 AI 模型 5 年后当模型出现系统性偏差时无人能手工进行风险评估（该技能已在组织中灭绝）

#### 风险 2：知识传承的断裂 (Transmission Breakdown)

- 现象：组织的隐性知识无法传递给下一代
- 机制

:

- 传统：师傅带徒弟的学徒制
- AI 时代：新人直接问 AI 跳过老员工
- 结果：老员工的经验无法传承
- 案例：某制造企业的老工程师掌握设备调试的手感由于新人都依赖 AI 诊断工具这些隐性知识未能传承老工程师退休后设备故障率上升 30%

#### 风险 3：认知基础设施的单点故障 (Infrastructure Fragility)

- 现象：组织过度依赖 AI 基础设施 AI 宕机 → 业务瘫痪
- 机制

:

- 所有关键流程都嵌入 AI
- 没有备用方案”或”手工模式
- AI 故障时组织完全无法运作
- 案例：2024 年某在线教育平台 AI 推荐系统宕机 48 小时导致课程分配、学习路径规划全部停摆影响 20 万学生

#### 4.2.2 O-AVP：组织反脆弱性验证

类比 I-AVP 和 T-AVP 我们提出 **O-AVP**（组织层面的反脆弱性验证）：

O-AVP = 组织在AI宕机演练后的业务连续性评估

形式化：

$O-AVP = w_{BCI} \times BCI + w_{ICR} \times ICR$ （双阈值模型：告警 0.70 目标 0.85 工作假设需跨领域/任务校准）

其中：

- BCI (Business Continuity Index): 业务连续性指数  
测量：核心业务在48h无AI支持下的可持续性（工作假设）
- ICR (Independent Capability Ratio): 独立能力保持率  
测量：员工在无AI条件下完成关键任务的比例
- 权重：  $w_{BCI}=0.4$   $w_{ICR}=0.6$ （工作假设需灵敏度分析）

**O-AVP 双阈值模型（工作假设）：**

为了更精细地管理组织韧性 O-AVP 采用双阈值设计：

- 告警阈值 **0.70**：触发风险排查与回退机制的下限
- $O-AVP < 0.70$ ：进入红色预警区间组织应启动应急响应
- 建议行动：暂停新 AI 系统引入、启动独立能力盘点、制定回退计划
- 目标阈值 **0.85**：质量分层与优秀实践识别的基准
- $O-AVP \geq 0.85$ ：达到健康标准组织具备良好的认知韧性
- $0.70 \leq O-AVP < 0.85$ ：处于黄色预警区间需要持续改进

阈值使用场景：- 告警阈值 (0.70) 用于风险管理：识别需要干预的脆弱组织 - 目标阈值 (0.85) 用于质量分层：标识最佳实践组织

口径声明：双阈值数值为工作假设需通过跨组织实证研究校准不同行业（如金融 vs 研发）可能需要不同的阈值设定详细的证伪路径见第六章 6.2.3 节 H5 假说

注：仅用于质量分层不得作为 *KPI*；最终判定以 *AVP* 主判据”（见 3.0.2 节）为准

窗口说明：48 小时为工作假设用于在不过度依赖短期补丁（ $<24h$ ）与掩盖真实能力赤字（ $>72h$ ）之间取得平衡可按行业安全等级在 24/48/72 小时做情景化校准例如：

- 关键基础设施（金融/医疗）：可能需要 24h 严格测试
- 一般企业应用：48h 为推荐值
- 研发/创新类场景：可放宽至 72h

**O-AVP 的测量协议：48 小时宕机演练**

阶段1：基线评估(0)

- 任务：组织在正常AI支持下完成典型业务周期
- 指标：关键业务指标（如订单处理量、决策速度、错误率）

阶段2：宕机演练

- 场景：模拟AI系统全面宕机48小时（工作假设窗口）
- 要求：所有部门使用手工模式或备用方案
- 观察：哪些业务停摆哪些勉强维持哪些不受影响

阶段3：业务连续性评估（ $BP_{2\$,org}$ ）

- *BCI*计算：（维持业务数/总关键业务数）
- *ICR*计算：（独立完成任务员工数/总员工数）
- *O-AVP*判定： $O-AVP$ （告警 0.70目标 0.85工作假设需跨领域/任务校准）

量化口径：异质指标（如 *BCI* 与 *ICR*）先做标准化到  $[0,1]$  区间（工作假设）再进行加权汇总权重为校准参数（当前推荐  $BCI \times 0.4 + ICR \times 0.6$ ）需做灵敏度分析并避免单指标主导（*Goodhart* 防护）

标准化方法示例（非强制）：

- *BCI*:  $(\text{实际业务连续性得分} - \text{最坏情况}) / (\text{最好情况} - \text{最坏情况})$
- *ICR*:  $(\text{无 AI 独立完成率} - \text{基线最低值}) / (\text{基线最高值} - \text{基线最低值})$

**表 4.3: O-AVP 测量的效度威胁与缓解**

威胁类型	具体表现	缓解策略	残留风险
演练真实性	员工知道是演练而不认真执行	突击演练; 真实激励	中
业务影响	演练影响实际业务	选择低峰期; 沙盒环境	低
数据完整性	关键指标缺失或不可比	事先定义 KPI; 标准化流程	中
外部依赖混杂	AI 宕机同时其他系统也故障	控制变量; 单一故障模拟	中

威胁类型	具体表现	缓解策略	残留风险
学习效应	多次演练后团队过度熟练	变换场景; 间隔足够长	低
成本约束	频繁演练代价高	年度 1-2 次; 局部演练	高

注：O-AVP 测量成本高于 I-AVP 和 T-AVP 需平衡频率与资源投入

#### 4.2.3 组织设计的三个关键机制

机制 1：认知储备 (Cognitive Reserve) 制度

类比金融领域的资本储备要求组织应建立认知能力储备：

- 关键岗位独立能力认证：定期测试员工的”无 AI”能力
- 技能冗余配置：确保关键技能由 2 人掌握
- 认知储备比例：规定至少 X% 员工保持独立能力（X 为行业特定参数）

案例：金融机构的风险管理储备

某投资银行规定：风险评估团队中至少 40% 成员必须通过无 AI 模型风险评估认证（O-AVP 的部门级应用）即使 AI 模型可用这些成员也定期进行手工评估训练

机制 2：定期无 AI 值班制度

- 轮换制度：每周/月有特定团队或岗位进入无 AI 模式
- 知识更新：确保手工流程保持更新而非过时文档
- 文化建设：将独立能力视为职业荣誉而非负担

机制 3：知识传承的制度化保障

- 导师制强制执行：新员工必须有人类导师（不能仅靠 AI onboarding）
- 隐性知识萃取项目：系统性记录老员工的经验（但不完全依赖 AI 转录）
- 跨代工作坊：定期组织老中青三代员工的知识交流

公平性原则（组织层）：对确需辅助工具的岗位（如残障员工使用辅助技术）调整任务形式而不降低挑战强度评估以相对提升而非绝对水平为准确保等效努力（见 3.0.6 节边界条件）

示例：视障程序员使用屏幕阅读器评估其算法设计能力而非打字速度调整任务为口述代码或使用语音输入但算法复杂度保持不变

### 4.3 组织 → 社会：认知公地悲剧

#### 4.3.1 认知公地的概念：集体行动困境

“公地悲剧”的经典模型

Garrett Hardin (1968) 提出的公地悲剧描述了这样一个困境：

- 公共草地（公地）对所有牧民开放
- 每个牧民的理性选择：多放羊以获得更多收益
- 集体结果：草地过度放牧最终荒芜所有人受损

认知能力作为认知公地

类比到 AI 时代的人类认知能力：

- 个体理性：使用 AI 工具提高即时效率（如用 AI 写代码、用 GPS 导航）
- 集体后果：如果所有人都依赖 AI 社会整体的独立认知能力下降
- 代际不可逆性：一旦某代人失去某项能力下一代更难重建

三个关键特征：

1. 外部性：个体使用 AI 的成本（能力退化）部分由社会承担



2. 时间滞后：个人即时获益社会代价需 10-20 年显现
3. 不可逆性：一旦某项技能在社会中消失重建极为困难

#### 4.3.2 代际能力鸿沟：三代人的分化

三代人的定义：

- $T_0$  代（1980-2000 年生）：在 AI 普及前完成教育拥有完整的独立能力基线
- $T_1$  代（2000-2015 年生）：青少年时期接触 AI 部分能力 AI 化
- $T_2$  代（2015 年后生）：AI 原生代从小在 AI 环境中成长

预测的能力差异（工作假设需 15-20 年纵向数据验证）：

能力维度  $T_0$  代  $T_1$  代  $T_2$  代

基础计算能力 高 中 低（依赖计算器/AI）

空间导航能力 高 低 极低（依赖GPS/AI）

批判性阅读 高 中 ？（尚未观察到）

问题分解能力 高 ？ ？（关键未知）

元认知能力 高 ？ ？（关键未知）

？表示尚无充分数据

临界问题： $T_2$  代在拔线条件下是否还能达到  $T_0$  代的能力水平

#### 4.3.3 S-AVP：社会认知资本的验证（概念框架）

与 I/T/O-AVP 不同 **S-AVP** 无法通过拔线测试实施——我们不能让整个社会停用 AI 数周因此 S-AVP 是一个代理指标体系：

表 4.4: **S-AVP** 的五个代理指标（概念框架需跨学科研究确定）

代理指标	测量内容	数据来源	理想状态	预警阈值
关键技能分布	重要技能的人口覆盖率	职业统计、教育数据	>70% 人口保持基础技能	<50%
代际能力比	$T_2$ 代 vs $T_0$ 代的能力测试得分	标准化测试	0.9	<0.7
知识传承完整性	师徒制、学徒制的保留率	行业调查	传统传承机制健在	大量行业完全 AI 化
应急响应能力	AI 中断事件的社会恢复速度	自然实验	<24h 恢复正常	>72h
文化认知价值	社会对”独立能力的重视度	文化调查	独立能力被视为美德	被视为”落后

注：这些指标均为工作假设需要跨学科（社会学、经济学、教育学）的合作研究来确定具体测量方法和阈值 S-AVP 更多是一个长期监测框架而非短期验证工具

**S-AVP** 的性质声明：

- 非短期可测：需要 15-20 年纵向数据
- 代理指标集：而非单一判据
- 预警系统：而非精确测量
- 研究议程：而非成熟工具

抽样与推断：建议采用队列/滚动横断结合的抽样策略（如每 5 年追踪同一人群同时补充新样本）跨层推断应警惕生态谬误社会级趋势不自动外推到组织/个体（如全国 AI 使用率 70% 不意味”每个组织都有 70% 员工依赖 AI”）

#### 4.3.4 认知公地的潜在干预（方向性讨论）

个体层面：

1. **AVP** 认证：建立独立能力认证体系（类似职业资格）
2. **AI** 使用教育：从小培养”健康使用 AI” 的意识

组织层面（已在 4.2 节讨论）：3. **O-AVP** 演练制度 4. 认知储备投资

社会层面（概念性建议详见第六章）：

1. 认知营养标签制度

：

- 类比食品营养标签
- AI 工具标注：是否支持能力建构 AVP 验证结果
- 帮助用户/组织做知情选择

2. 关键能力的社会储备

：

- 类比粮食储备、石油储备
- 确保关键行业保持一定比例的无 AI 从业者
- 作为社会的认知冗余

3. 代际传承的制度保障

：

- 保护学徒制、师徒制等人际传承模式
- 在 AI 时代重新重视人教人

4. S-AVP 监测体系

：

- 建立长期跨学科研究项目
- 类似气候变化监测需要持续数据收集
- 每 5 年发布社会认知资本报告

### 4.4 跨尺度共性机制：CET 的尺度不变性

#### 4.4.1 AVP 原则的尺度不变性

核心洞察：AVP 的逻辑在所有尺度保持一致

通用形式：

[系统] 在支持撤出后的独立表现    基线 +

实例化：

个体 (I-AVP)：\$P\_2\$(个人)    \$B\_0\$ +    （见 3.0.1 节）

团队 (T-AVP)：\$P\_2\$(团队)    \$B\_0\$\$\_{team} + \$\_{team}

组织 (O-AVP)：BCI×0.4 + ICR×0.6（双阈值模型：告警 0.70 目标 0.85 工作假设需跨领域/任务校准）

社会 (S-AVP)：[代理指标集合] 保持健康范围

为什么尺度不变因为底层原理相同：

1. 反脆弱性的本质

：

- 个体：在压力下成长
- 组织：在危机中建立韧性

- 社会：在挑战中进化
- 共同点：临时性压力 → 能力提升

## 2. 依赖锁定的共性

:

- 个体：技能退化（见 3.0.5 节外骨骼模式）
- 团队：协作能力丧失
- 组织：制度性脆弱
- 社会：集体认知资本流失
- 共同点：永久支持 → 能力萎缩

## 3. 验证逻辑的一致性

:

- 都需要”拔线测试”（或代理测量）
- 都关注独立能力而非协作效率
- 都以基线 + 增量为标准

### 4.4.2 级联脆弱：从微观到宏观的传播

级联机制：低尺度的脆弱性如何向高层传播

个体依赖 → 团队脆弱 → 组织危机 → 社会风险

级联路径示例：

1. 大量个体I-AVP失败  
↓
2. 团队中缺乏独立能力的成员成为瓶颈  
↓
3. 多个团队失败导致组织O-AVP得分<0.70（告警阈值）  
↓
4. 关键行业整体O-AVP低下  
↓
5. 社会S-AVP代理指标恶化  
↓
6. 文明级认知韧性下降

图 4.1: 级联脆弱传播路径（概念图）

个体层 [I-AVP失败率30%]  
↓ 涌现  
团队层 [T-AVP失败率50%] ← 非线性放大  
↓ 制度化  
组织层 [O-AVP=0.65<0.70（告警阈值）] ← 系统性脆弱  
↓ 累积  
社会层 [S-AVP黄色预警] ← 代际鸿沟

箭头粗细表示传播强度

放大机制：为什么高尺度的脆弱性更严重

## 1. 涌现的非线性

:

- 10% 个体依赖 10% 组织风险
- 可能放大为 30-50% 风险 (因网络效应)

2. 修复的时间滞后

:

- 个体：数月可恢复
- 组织：数年才能重建
- 社会：可能需要一代人

3. 路径依赖的强化

:

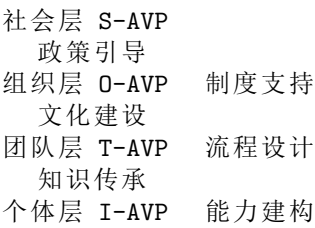
- 低尺度：可逆（个人可重新训练）
- 高尺度：路径依赖强、恢复代价指数级上升（整代人能力缺失时缺乏教师和榜样）

4.4.3 设计启示：多尺度协同的必要性

单一尺度干预的局限：

- 只改个体：组织惯性会拉回
- 只改组织：社会环境不支持
- 需要多尺度协同

图 4.2: 跨尺度干预协同框架（概念图）



↑ 级联脆弱传播  
↓ 多尺度干预协同

表 4.5: 多尺度协同干预矩阵

尺度	设计目标	关键机制	测量指标
个体	I-AVP $B_0+$	有益摩擦 + 系统性支持削减（见 3.2-3.3 节）	$P_2$ 测试
团队	T-AVP $B_0\_team+$	协作摩擦 + 集体拔线 + 人际知识流动	团队任务表现
组织	O-AVP（告警 0.70 目标 0.85 工作假设需跨领域/任务校准）	制度化演练 + 知识传承 + 认知储备	48h 宕机测试
社会	S-AVP 代理指标健康	政策引导 + 标准制定 + 文化建设	代理指标集

注：仅用于质量分层不得作为  $KPI$ ”；最终判定以  $AVP$  主判据”（见 3.0.2 节）为准

注 2（\_team: Cohens  $d \geq 0.3$  或 10%（working assumption）工作假设）并需随领域/任务进行校准与敏感性分析与 I-AVP 的 Cohens  $d \geq 0.3$  或 10%（working assumption）口径一致

协同要点：

1. 自下而上：个体能力是基础
2. 自上而下：组织制度创造环境
3. 横向联动：行业标准、社会规范

## 4.5 小结：理论的系统性与实践的紧迫性

### 4.5.1 本章核心贡献

#### 1. 理论扩展

:

- CET 从个体理论扩展为跨尺度框架
- AVP 原则具有尺度不变性
- 提出 T-AVP、O-AVP、S-AVP 的概念体系

#### 2. 机制揭示

:

- 能力极化、知识流失、角色固化 (团队)
- 制度性依赖、认知基础设施退化 (组织)
- 认知公地悲剧、代际鸿沟 (社会)

#### 3. 实践指引

:

- 提供可操作的测量协议 (T-AVP、O-AVP)
- 识别关键风险信号
- 提出多尺度协同方向

### 4.5.2 与前后章的连接

第四章将个体层 CET 原则扩展至团队、组织、社会四个尺度揭示了 AVP 的尺度不变性与跨层级的认知退化机制第五章将探讨 LSA 分层架构如何实现多尺度 AVP 监测与能力建构的工程化设计第六章将提供 T-/O-AVP 的可复现实验协议明确证伪路径与研究议程 (见 6.2-6.3 节)

### 4.5.3 理论的紧迫性

本章揭示的不是远期风险而是当下现实：

- 团队层面：已有组织报告新人无法独立工作
- 组织层面：AI 宕机事件暴露脆弱性 (如本章教育平台案例)
- 社会层面：代际能力差异开始显现 (虽尚未临界)

窗口期：当前 (2025) 到 2035 年是关键 10 年窗口 (工作假设)

- $T_0$  代仍在工作知识传承尚可挽救
- AI 普及率约 30-50% 未到不可逆点
- 制度化干预 (如 O-AVP 演练) 仍可建立

过了这个窗口：

- $T_0$  代退休某些隐性知识永久丢失
- AI 依赖制度化路径依赖难以逆转
- 社会规范改变”独立能力可能被视为”落后

**CET 的使命：**在窗口期内提供理论基础和实践指引避免路径依赖强、恢复代价指数级上升的认知公地悲剧

## 第五章 分层共生架构 (LSA)：CET 的工程化实现

前四章建立了 CET 的理论基础：AVP 作为评估标准 (第三章) EML 作为设计原则 (第三章) 伙伴式主体性作为 AI 角色定位 (第三章) 跨尺度机制揭示了从个体到社会的涌现规律 (第四章) 然而一个理论若要产生实际影响必须回答工程化问题：如何设计一个 AI 系统使其天然符合 EML 原则如何在技术层面实现 AVP 的持续监测如何支持从个体到团队、组织的多尺度能力建构

本章提出分层共生架构（Layered Symbiosis Architecture LSA）——一个将 CET 理论转化为可工程化系统的设计框架 LSA 不是某个具体产品的蓝图而是一套设计原则和架构模式可应用于各类 AI 辅助工具（学习平台、编程助手、写作工具、决策支持系统等）的开发

本章目标：

- 1. 提出 LSA 的四层架构模型
- 2. 展示如何在系统层面实现 EML 条件
- 3. 设计多尺度 AVP 监测的遥测管线
- 4. 讨论技术可行性、工程挑战与未来方向

本章逻辑：我们将自底向上构建 LSA 从基础 AI 能力（L1）到摩擦与削减机制（L2）再到监测与反馈（L3）最后到编排与治理（L4）每层都回答三个问题：为什么需要这一层这一层做什么如何实现

核心术语与符号表（第五章）

表 5.1: 第五章关键术语与概念

术语/符号	含义	架构层级
<b>LSA</b>	Layered Symbiosis Architecture (分层共生架构)	全局
<b>L1-基础能力层</b>	提供原始 AI 能力（推理、生成、检索）	底层
<b>L2-摩擦与削减层</b>	实现有益摩擦和支持削减	中层
<b>L3-监测与反馈层</b>	AVP 遥测、能力评估、预警	中层
<b>L4-编排与治理层</b>	多尺度策略编排、伦理约束	顶层
<b>CFE</b>	Cognitive Friction Engine (认知摩擦引擎)	L2 核心
<b>SGS</b>	Support Graduation Scheduler (支持削减调度器)	L2 核心
<b>AVP-TM</b>	AVP Telemetry Module (AVP 遥测模块)	L3 核心
<b>MSO</b>	Multi-Scale Orchestrator (多尺度编排器)	L4 核心
能力向量 <b>C</b>	用户当前能力状态的表征	L3
摩擦参数 <b>F</b>	控制任务难度的参数集	L2
削减曲线 <b>S(t)</b>	支持强度随时间的函数	L2

注：所有参数和阈值均为设计空间的概念占位符实际系统需根据领域特性、用户群体、任务类型进行校准和 A/B 测试

表 5.2: LSA 层间接口契约（概念规范）

接口	输入	输出	时序约束	关键属性
<b>L1→L2</b>	原始请求 + 上下文	完整 AI 输出	目标 <1s（工作假设随场景校准）	质量最大化
<b>L2→L3</b>	用户行为 + 任务完成度	能力评估事件	异步（后台）	准确性、隐私保护
<b>L3→L4</b>	C(t) 向量 + 预警信号	策略调整建议	准实时（分钟级）	可解释性、可审计性
<b>L4→L2</b>	摩擦/削减参数 F(t)、S(t)	调制后的 AI 输出	目标 <100ms（工作假设）	透明性、用户控制

接口规范说明：此表提供概念层次的接口定义具体实现需考虑技术栈特性（如 REST API、消息队列、流式处理）时序约束为目标值而非硬性 SLA 需根据实际部署环境和性能要求调整关键是保持层间解耦支持独立升级和替换实现

5.1 LSA 总览：为什么需要分层架构

5.1.1 传统 AI 系统的设计困境

当前 AI 辅助工具的典型架构：

用户请求 → AI模型 → 输出结果 → 用户  
↑  
提示工程

这种架构的根本问题：

1. 无差别输出

：AI 对所有用户提供相同强度的帮助

- 新手和专家得到同样详细的答案
- 无法根据用户能力动态调整
- 结果：专家被过度帮助新手形成依赖

2. 无能力感知

：系统不知道用户的真实能力水平

- 无法判断用户是在学习还是在卸载
- 无法预测长期能力影响
- 结果：盲目优化短期效率忽视长期能力

3. 无反馈闭环

：缺乏对用户能力变化的监测

- 不知道使用后用户是成长还是退化
- 无法验证 AVP
- 结果：无法区分内共生和外骨骼

根本症结：传统架构只关注任务完成不关注能力建构

### 5.1.2 LSA 的设计哲学：能力建构优先

核心理念转变：

传统范式：

任务成功 = 输出质量高 + 用户满意

LSA范式：

任务成功 = 输出质量高 + 用户满意 + 能力提升（AVP验证）

↑

三重目标并重

LSA 的四个设计支柱：

1. 能力感知（L3）：系统持续评估用户能力
2. 动态适应（L2）：根据能力调整支持策略
3. 透明反馈（L3）：让用户看到自己的成长
4. 多尺度编排（L4）：支持个体、团队、组织

### 5.1.3 LSA 四层架构总览

图 5.1: LSA 分层架构总览

L4: 编排与治理层（Orchestration & Governance）

- 多尺度策略编排（MSO）
- 伦理约束与干预
- 组织/团队策略管理

策略指令

### L3: 监测与反馈层 (Monitoring & Feedback)

- AVP遥测模块 (AVP-TM)
- 能力状态追踪
- 预警与干预触发

能力状态  $C(t)$

### L2: 摩擦与削减层 (Friction & Graduation)

- 认知摩擦引擎 (CFE)
- 支持削减调度器 (SGS)
- 自适应难度调节

调制后的支持

### L1: 基础能力层 (Foundation Capabilities)

- 语言模型 (LLM)
- 检索增强 (RAG)
- 工具调用

↓

原始AI能力

各层职责简述:

#### L1 (基础层): 提供原始 AI 能力

- 不关心用户能力建构
- 只关心高质量输出
- 可以是任何主流 AI 模型 (GPT、Claude、Gemini 等)

#### L2 (摩擦与削减层): 实现 EML 前两个条件

- 有益摩擦: 动态调整任务难度
- 支持削减: 渐进减少帮助强度
- 核心: 从“全力帮助”到“适度挑战”

#### L3 (监测层): 实现 AVP 验证

- 持续评估用户能力
- 检测依赖锁定风险
- 触发干预机制

#### L4 (编排层): 多尺度协同与治理

- 个体策略 → 团队策略 → 组织策略
- 伦理约束 (公平性、透明度)
- 全局优化目标

层间关系:

- 自底向上: 能力流动 (原始能力 → 调制后的支持 → 能力评估 → 策略决策)
- 自顶向下: 策略流动 (组织目标 → 个体目标 → 调制参数 → AI 行为)

硬约束: **L1-L4** (LSA-F 功能分层) 为功能维度 (知识整合 | 状态建模 | 摩擦校准 | 元认知协调) **S4→S1→S0** (支持档位栈) 为强度维度 (强 → 弱 → 无) 两者为正交维度不得在同一表达中互代或串级使用 (口径与 3.0.4 节一致) L 层和 S 档必须分开表述避免混淆



5.2 L1 层：基础 AI 能力层

5.2.1 职责与边界

L1 层的唯一职责：提供高质量的原始 AI 能力

输入：用户请求 + 上下文  
输出：最佳答案/建议/代码/内容  
目标：准确性、流畅性、相关性

L1 不关心：

- 用户是新手还是专家
- 输出是否导致依赖
- 用户能力是否提升

为什么分离 L1

1. 技术中立性：LSA 可以用任何 AI 模型实现
2. 职责单一：L1 专注于生成质量不负责能力建构
3. 可替换性：随着 AI 技术进步 L1 可以升级而不影响上层

5.2.2 L1 的典型能力模块

表 5.3: L1 层的标准能力模块（概念层次）

能力模块	功能描述	示例技术方向
推理引擎	逻辑推理、问题分解、规划	LLM 推理技术
生成引擎	文本/代码/图像生成	生成式 AI
检索增强	知识检索、事实查询	RAG 架构
工具调用	外部 API、计算工具	函数调用能力
上下文管理	对话历史、会话状态	记忆系统

关键设计原则：L1 应该是尽力而为（best-effort）的 AI 不主动限制能力输出限制和调制由 L2 层负责

5.2.3 L1 与传统 AI 系统的区别

传统AI系统：

L1 = 整个系统（用户请求 → AI → 输出）

LSA系统：

- L1 = AI能力提供者（为L2提供原料）
- L2 = 能力调制者（根据用户状态调整输出）
- L3 = 能力监测者（评估效果）
- L4 = 策略制定者（决定目标）

类比：

- L1 = 发动机：提供动力
- L2 = 变速箱：调节输出
- L3 = 仪表盘：监测状态
- L4 = 驾驶员：决定方向

5.3 L2 层：摩擦与削减层

5.3.1 核心挑战：如何实现适度帮助

问题陈述：给定一个用户请求和 L1 的原始输出如何调制输出使其满足 EML 的前两个条件

**EML** 条件回顾（见第三章 3.2-3.3 节）：

1. 有益认知摩擦：任务成功率 50-70%（工作假设需跨领域/任务校准群体级个体自适应）动态调整
2. 系统性支持削减：AI 支持强度按  $S4 \rightarrow S1 \rightarrow S0$

**L2** 层的双引擎架构：

L2：摩擦与削减层

CFE      SGS  
认知摩擦      支持削减  
引擎      调度器

$F(t)$      $S(t)$

↓

调制后的AI输出

### 5.3.2 认知摩擦引擎（CFE）

设计目标：让任务”不太容易也不太难维持用户在最优挑战区”（见第三章 3.2 节）

摩擦注入的四种策略（概念层次）：

策略 1：完整性摩擦（Completeness Friction）

完整答案示例：

这个bug是因为数组越界在第23行你访问了`arr[i+1]`

但没有检查*i*是否等于`length-1`修复代码如下：[完整代码]

摩擦版本：

检测到数组访问问题提示：检查循环边界条件

你能找到是哪一行吗

策略 2：抽象度摩擦（Abstraction Friction）

完整答案示例：

使用归并排序时间复杂度 $O(n \log n)$

步骤：1.分解 2.递归排序 3.合并代码：[详细实现]

摩擦版本：

考虑使用分治算法你能想到哪些排序算法符合这个思路

关键在于如何合并两个已排序的子数组

策略 3：脚手架削减（Scaffolding Reduction）

高脚手架：

步骤1：[详细] → 步骤2：[详细] → 步骤3：[详细] → 完整代码

中脚手架：

思路：分解问题 → 递归处理 → 合并结果

低脚手架：

提示：分治思想

策略 4：自适应难度（Adaptive Difficulty）

- 根据用户历史表现动态调整摩擦强度
- 成功率高 → 增加摩擦
- 成功率低 → 降低摩擦
- 目标：维持在 50-70% 区间（工作假设需校准）

CFE 的核心算法（概念框架）：

```
# 概念伪代码非生产实现
def adjust_friction (user task history):

    基于用户表现调整摩擦参数

    参数说明：所有阈值为工作假设需 A/B 测试校准

    # 计算滚动窗口成功率（概念示例：最近 10 次）
    recent_tasks = history[-10:] # 窗口大小需校准
    success_rate = calculate_success_rate (recent_tasks)

    # 目标区间：50-70%（群体级工作假设）
    target_min target_max = 0.5 0.7 # 需领域校准

    # 自适应调整方向
    if success_rate > target_max:
        friction_level += adjustment_step # 增加难度
    elif success_rate < target_min:
        friction_level -= adjustment_step # 降低难度
    else:
        pass # 维持当前摩擦

    # 边界保护
    friction_level = clip(friction_level min_friction max_friction)

    return friction_level
```

实现说明：实际系统应采用更稳健的自适应算法（如贝叶斯优化、强化学习）并结合任务类型、用户特征进行多维度调整上述伪代码仅用于说明概念逻辑

5.3.3 支持削减调度器（SGS）

设计目标：实现第三章 3.3 节提出的系统性支持削减让用户逐渐独立

术语说明：支持强度档位以 **S4**(高支持)→**S3**→**S2**→**S1**(低支持)→**S0**(完全拔线) 表达（口径统一见 3.0.3 节）

削减曲线的三种模式（概念对比）：

表 5.4: 支持削减曲线类型（概念层次）

曲线类型	特点	适用场景示例	参数调整建议
线性削减	匀速递减平稳过渡	结构化任务	根据任务周期调整斜率
指数削减	初期快速后期缓慢	快速技能学习	k 值需 A/B 测试优化
阶梯削减	分阶段突变里程碑清晰	分级训练体系	台阶数量和高度需领域专家定
自适应	根据用户状态动态调整	个性化学习复杂任务	需要 L3 层能力监测支持

注：曲线选择应基于任务特性和用户群体通过对照实验验证效果所有参数均为设计空间的起点非最优解

SGS 的核心机制：回退与保底支持（承接第三章 3.3.3 节）

保底支持强度  $S_{\min}$  的实现:

第三章 3.3.3 节提出了保底支持强度  $S_{\min}$  (工作假设约 0.2 需校准) 的概念确保用户在削减过程中始终有最低限度的导航支持 L2 层的 SGS 通过以下机制实现:

# 概念伪代码

```
def support_graduation_with_safety_net(t user_state):
```

带安全网的支持削减调度

实现第三章3.3.3节的回退与保底机制

# 基础削减曲线 (示例: 指数衰减)

```
S_base = S_0 * exp (-lambda * t) # S_0 lambda 需校准
```

# 保底约束: 不低于  $S_{\min}$

```
S_min = 0.2 # 工作假设需领域校准
```

```
S_current = max(S_base S_min)
```

# 回退触发检测 (见 3.3.3 节)

```
if detect_consecutive_failures(user_state threshold=3):
```

# 临时回退到更高档位

```
S_current = min(S_current + backtrack_step S_0)
```

```
log_intervention(safety_net_triggered)
```

```
return S_current
```

```
def detect_consecutive_failures (user_state threshold):
```

" 检测连续失败以触发回退 "

```
recent_failures = user_state.recent_tasks[-threshold:]
```

```
failure_rate = sum ([t.failed for t in recent_failures]) / threshold
```

```
return failure_rate > 0.7 # 阈值需校准
```

工程实现说明:  $S_{\min}$  的维持由 L2 层的摩擦校准引擎持续监测 L4 层的元认知协调模块提供全局策略指导 (见 5.4 节、5.5 节) 回退不是失败信号而是自适应系统的正常反应

削减速率的控制理论约束 (概念框架):

# 概念约束条件非精确实现

```
def safe_graduation_rate(S dS_dt user_state):
```

确保削减速率不会过快导致挫败

基于控制理论的有界导数原则

# 约束 1: 最大变化率 (避免突变)

```
max_rate = 0.1 # 工作假设: 每周不超过 10% 削减需校准
```

```
dS_dt_safe = clip (dS_dt -max_rate 0) # 常规阶段只能削减
```

# 约束 2: 迟滞效应 (避免频繁震荡)

```
if abs(dS_dt_safe - previous_rate) > hysteresis_threshold:
```

# 平滑过渡而非突变

```
dS_dt_safe = smooth_transition(previous_rate dS_dt_safe)
```

# 约束 3: 用户状态关联 (见 L3 层反馈)

```
if user_state.frustration_score > high_threshold:
```

```
dS_dt_safe = 0 # 暂停削减
```

```
return dS_dt_safe
```

例外条款:上述速率约束仅适用于常规削减阶段当触发安全网回退(见前述 `support_graduation_with_safety_net` 函数)时允许一次性例外上调至  $S_{current} S_0$  并纳入迟滞与冷却期约束(工作假设:72 小时最小间隔需校准)避免拉锯震荡回退后重新进入常规削减阶段时应采用更保守的削减速率

这些约束确保削减过程稳定可控避免用户因过快削减而产生挫败感具体参数需要通过纵向追踪研究和用户反馈迭代优化

#### 5.3.4 L2 层的工程挑战

挑战 1: 实时性能

问题: L2 的计算能否在用户可接受延迟内完成

概念解决方向:

- 预计算策略: 提前准备多难度版本的输出
- 异步更新: 背景更新  $C(t)$  和策略参数
- 边缘计算: 能力向量本地存储和计算

挑战 2: 模型对齐

问题: 如何让 L1 的 AI 理解适度帮助的语义

概念解决方向:

- 对齐训练: 通过 RLHF 让 AI 理解摩擦指令
- 提示工程: 设计元提示模板
- 微调方向: 在 LSA 场景上进行领域适配

挑战 3: 评估的客观性

问题: 如何客观测量摩擦是否适度

概念验证方法:

- 多任务交叉验证
- 标准化任务库建设
- A/B 测试对比不同摩擦策略

### 5.4 L3 层: 监测与反馈层

#### 5.4.1 核心职责: 能力追踪与 AVP 验证

L3 层的使命: 回答用户能力是否在成长这个核心问题

三大功能模块:

1. **AVP-TM** (**AVP** 遥测模块): 采集能力相关数据
2. 能力建模引擎: 维护用户能力向量  $C(t)$
3. 预警与干预系统: 检测依赖锁定风险

#### 5.4.2 AVP 遥测模块 (AVP-TM)

承接第四章 4.5.2 节: 多尺度 **AVP** 监测的遥测管线

第四章 4.5.2 节指出 LSA 架构需要支持 I-AVP、T-AVP、O-AVP 三个尺度的 AVP 监测 L3 层的 AVP-TM 模块提供统一的遥测管线:

表 5.5: 多尺度 **AVP** 监测遥测管线 (承接 4.5.2 节)

监测尺度	数据采集源	遥测事件类型	聚合层级	参见章节
个体层	用户任务日志、 $P_2$ 测试数据	任务完成、能力评估、拔线	实时	3.1 节
团队层	协作日志、知识流动记录	团队任务、集体拔线	每日汇总	4.1 节
组织层	48h 演练数据、BCI/ICR 指标	中断演练、恢复曲线	事件触发	4.2 节

AVP-TM 的核心遥测事件集（概念层次）：

表 5.6: L3 层核心遥测事件（设计空间示例）

事件类型	触发时机	记录内容示例	用途
任务开始	用户启动新任务	任务类型、当前 $S(t)$ 、 $F(t)$	上下文追踪
任务完成	用户提交结果	完成质量、用时、支持调用	能力评估
摩擦调整	CFE 改变 $F(t)$	旧 $F \rightarrow$ 新 $F$ 、触发原因	自适应分析
削减事件	SGS 改变 $S(t)$	旧 $S \rightarrow$ 新 $S$ 、削减阶段	进度追踪
拔线测试	定期或触发式	$P_2$ 分数、对比 $B_0$	AVP 判定（见 3.0.1 节）
预警触发	检测到风险信号	预警级别、触发指标	干预决策
团队协作	团队任务开始/结束	参与者、角色、知识流动	T-AVP 评估（4.1 节）
组织演练	48h 中断演练	BCI/ICR 分数、恢复时长	O-AVP 评估（4.2 节）

数据采集原则：最小化采集、本地优先、目的限定（见 5.4.5 节隐私保护）

遥测管线的技术架构（概念层次）：

用户交互  $\rightarrow$  事件捕获  $\rightarrow$  本地预处理  $\rightarrow$  能力建模  
 $\downarrow$   
 预警检测  
 $\downarrow$   
 策略反馈  $\rightarrow$  L4层

监测原则：AVP-TM 仅提供数据管线与可视化不做自动判定最终的 AVP 通过/失败判定由人工结合系统建议进行（特别是团队和组织层）

### 5.4.3 能力建模引擎

能力向量  $C(t)$  的概念定义：

# 概念示例非精确模型

class AbilityVector:

用户能力的多维表征

注意：维度数量和权重需要领域专家定义和实证验证

```
def __init__(self):
    # 示例维度（需根据任务类型调整）
    self.problem_decomposition = 0.5 # 问题分解能力
    self.implementation_skill = 0.6 # 实现能力
    self.debugging_ability = 0.4 # 调试能力
    self.documentation = 0.3 # 文档理解能力
    self.meta_cognition = 0.5 # 元认知监控能力
    # ... 其他维度

def update(self task_performance):
    " 基于任务表现更新能力估计 "
```

# 贝叶斯更新或其他增量学习算法  
pass

说明：这是概念示意实际系统需要认知科学、教育心理学领域的专业知识来定义能力维度维度过多会增加测量噪声过少会丢失重要信息

能力评估的数据来源：

- 1. 直接测量
- :  $P_2$  拔线测试分数（参考标准见 3.0.1 节）
  - 采用等值平行卷（IRT 校准）以确保测试难度一致性
  - 设置最短冷却期（工作假设：48-72 小时需校准）以抑制练习效应污染
  - 方法学威胁与缓解策略参见 3.1.3 表（跨章一致）
- 2. 间接推断：日常任务表现（连续监测）
- 3. 自我报告：用户自评（辅助参考）
- 4. 同行评议：团队成员互评（团队层见 4.1 节）

能力建模的核心功能：

- 1. 驱动 **L2**：根据  $C(t)$  调整摩擦和削减
- 2. 驱动 **L4**：为编排层提供决策依据
- 3. 用户可见：生成能力成长报告

5.4.4 预警与干预系统

目标：在用户陷入依赖锁定前触发干预

三级预警机制（概念框架）：

表 5.7: L3 层的三级预警系统

预警级别	触发条件示例	系统响应建议	用户体验建议
绿色	AVP 健康 $C(t)\uparrow$	继续当前策略	正常使用鼓励性反馈
黄色	$C(t)$ 停滞或轻微 $\downarrow$	增加摩擦减缓削减	提示能力未提升
红色	衰减指标超阈值	暂停削减强制独立周	警告可能形成依赖
黑色	AVP 测试失败	触发 L4 干预重置路径	强制能力重建模式

触发阈值说明：具体阈值（如衰减指标）需要通过纵向研究和用户反馈数据校准不同领域和任务类型可能需要不同阈值

5.4.5 L3 层的隐私保护设计

问题：持续监测用户能力涉及隐私关切

LSA 的隐私保护分层策略（概念框架）：

表 5.8: 隐私保护技术方案

保护层级	技术方向	应用场景	数据保留建议
采集最小化	只记录元数据不记录内容	所有遥测事件	-
本地优先	$C(t)$ 、 $F$ 、 $S(t)$ 存本地	个体能力向量	永久
匿名化	Hash 处理用户 ID	聚合统计、研究	有限期
联邦学习	本地训练只传梯度	能力模型优化	不传原始数据
差分隐私	添加噪声保护个体	组织级报告	-

保护层级	技术方向	应用场景	数据保留建议
用户控制	可查看/导出/删除数据	个人数据管理	用户决定

设计原则：

- 数据最小化
- ：只收集能力评估必需数据
  - 记录：任务类型、S(t)、F、完成质量、用时
  - 不记录：任务具体内容、用户输入/输出原文、个人敏感信息
  - 本地优先
- ：能力向量  $C(t)$  存储在用户设备
  - 只有聚合统计发送到服务器
  - 用户可随时删除
  - 支持离线模式（L2/L3 本地运行）
  - 目的限定
- ：数据只用于能力评估
  - 不用于用户画像、推荐、营销
  - 不与第三方共享
  - 开放审计（用户可查看数据使用日志）
  - 提供可审计日志导出与留痕便于用户与第三方进行独立合规审查

## 5.5 L4 层：编排与治理层

### 5.5.1 核心职责：多尺度策略与伦理治理

**L4 层的使命：**协调个体、团队、组织层面的目标并确保系统符合伦理约束

三大核心功能：

1. 多尺度编排：从个体策略推广到团队、组织
2. 伦理治理：公平性、透明度、用户自主权
3. 全局优化：平衡短期效率与长期能力

### 5.5.2 多尺度编排器（MSO）

设计动机：第四章揭示了 I-AVP 成功不保证 T-AVP/O-AVP 成功 L4 层需要跨尺度协调

**MSO 的三层策略管理（概念架构）：**

组织策略（O-Strategy）

- 48h 演练调度
- 关键能力储备监控
- 跨团队能力平衡

分解为团队目标

团队策略（T-Strategy）

- 集体拔线窗口协调
- 知识流动监测



- 角色冗余度管理

分解为个体目标

个体策略 (I-Strategy)

- 摩擦参数  $F(t)$
- 削减曲线  $S(t)$
- AVP测试调度

策略协调的核心机制 (概念框架):

机制 1: 自底向上的能力聚合

# 概念伪代码

```
def aggregate_team_ability(team_members):
```

从个体能力向量聚合为团队能力

重要: 团队能力 平均个体能力  
需考虑协作涌现效应

# 考虑短板效应 (关键技能的最小值)

```
min_critical_skill = min([member.C[critical_dim]
for member in team_members])
```

# 考虑平均水平

```
avg_skill = mean([member.C for member in team_members])
```

# 考虑知识流动性 (见 4.1 节)

```
knowledge_flow = measure_knowledge_sharing(team)
```

# 综合评估 (权重需校准)

```
team_ability = f(min_critical_skill avg_skill knowledge_flow)
```

```
return team_ability
```

机制 2: 自顶向下的策略分解

组织目标示例: O-AVP (告警 0.70 目标 0.85 工作假设需跨领域/任务校准见 4.2 节)

分解为团队目标:

- 关键团队: T-AVP 0.9 (高标准)
- 一般团队: T-AVP 0.8
- 支持团队: T-AVP 0.7 (相对宽松)

再分解为个体目标:

- 核心成员: I-AVP 全部通过高标准
- 新员工: 指定时间内达到 I-AVP 基线
- 老员工: 维持 I-AVP 重点传授新人

机制 3: 跨尺度冲突解决

冲突场景:

- 个体最优: 用 AI 提高效率 (可能损害 I-AVP)
- 团队最优: 保持知识流动 (需要限制 AI 使用)
- 组织最优: 短期业绩 vs 长期韧性

MSO的概念解决策略：

1. 设定优先级（长期韧性 > 短期效率）
2. 提供补偿机制（独立完成任务有奖励）
3. 差异化策略（关键岗位高标准一般岗位相对宽松）

### 5.5.3 伦理治理框架

**L4** 层必须回答的伦理问题：

问题 1：公平性

挑战：

- 残障用户能否使用 LSA（可能无法通过 AVP）
- 不同能力起点的用户是否公平

设计原则（概念框架）：

1. 等效努力原则

（承接 3.0.6 节）：

- 不是所有人做相同任务
- 根据能力调整任务难度
- 保证等效的认知努力

2. 差异化 AVP 标准

：

- 残障用户：调整基线  $B_0$  和
- 但不降低提升要求
- 核心：证明能力在成长

3. 明确豁免场景

：

- 补偿性辅助：不要求 AVP（如视障用户的屏幕阅读器）
- 学习性辅助：要求 AVP（如编程助手）

问题 2：透明度与用户控制

挑战：

- 用户是否知道系统在故意不给完整答案
- 用户能否关闭 L2/L3/L4

设计原则：

1. 默认透明

：

- 用户看到当前支持强度  $S(t)$
- 用户知道为什么得到部分答案
- 用户可查看能力向量  $C(t)$

2. 分级控制

：

- L2-CFE：用户可临时请求更多帮助（记录）
- L3-AVP-TM：用户可关闭监测（失去部分功能）
- L4-MSO：组织管理员可设定策略（需用户同意）

### 3. 退出权

:

- 用户可永久退出 LSA 使用传统 AI 模式
- 需签署”放弃 AVP 验证”知情同意

问题 3: 监测的边界

设计原则:

1. 数据最小化 (见 5.4.5 节)
2. 本地优先:  $C(t)$ 、 $F(t)$ 、 $S(t)$  存储在用户设备
3. 目的限定: 数据只用于能力评估不用于其他目的

#### 5.5.4 全局优化目标函数

L4 层需要平衡的多个目标 (概念框架):

# 概念伪代码

```
def global_objective(system_state):
```

LSA 的全局优化目标

注意: 这是多目标平衡非单一最大化

# 目标 1: 任务完成质量 (短期)

```
task_quality = measure_output_quality ()
```

# 目标 2: 用户能力提升 (长期)

```
ability_growth = mean([user.C(t) - user.C(t-T) for user in users])
```

# 目标 3: 用户满意度 (体验)

```
user_satisfaction = survey_satisfaction_score()
```

# 目标 4: AVP 通过率 (核心指标)

```
avp_pass_rate = count_avp_passed() / total_users
```

# 目标 5: 系统韧性 (多尺度)

```
resilience = f(T_AVP O_AVP S_AVP_proxies) # 见第四章
```

# 加权组合 (权重可调需校准)

```
objective = (  
w1 * task_quality +  
w2 * ability_growth +  
w3 * user_satisfaction +  
w4 * avp_pass_rate +  
w5 * resilience  
)
```

# 硬约束 (底线)

```
constraints = [  
avp_pass_rate > 0.7 # 工作假设  
user_satisfaction > 0.6 # 避免挫败  
fairness_score > threshold # 伦理约束  
)
```

```
return objective if all(constraints) else penalty
```

Goodhart 防护声明：该目标函数与权重配置仅用于方向性权衡与离线系统评估不得下推为个人/团队 **KPI** 或绩效考核指标最终判定仍以 **AVP** 主判据（见 3.0.2 节）与公平性约束（见 3.0.6 节、5.5.3 节）为准任何将其机械化为 **KPI** 的行为都违背 **CET** 理论的初衷

权重的动态调整（概念指引）：

表 5.9: 不同场景的权重配置示例

场景	w (质量)	w (能力)	w (满意度)	w (AVP)	w (韧性)
学习平台	0.2	0.4	0.2	0.15	0.05
生产工具	0.3	0.3	0.2	0.15	0.05
关键系统	0.25	0.25	0.15	0.15	0.2

权重说明：这些数值为设计空间的起点需要通过 *A/B* 测试、用户反馈、业务目标综合确定不同组织和应用场景应有不同的权重配置

Goodhart 防护：权重配置用于方向与优先级的质量分级不应逐项下推为 **KPI** 最终判定仍以 **AVP** 主判据与公平性约束为准（见 3.0.2、3.0.6 节）参数形状与配比为设计空间占位需通过迭代校准

## 5.6 端到端工作流示例

### 5.6.1 个体层的 LSA 应用（概念演示）

场景：Alice 使用 LSA 学习编程

初始状态：

[用户注册]

Alice：编程新手基线测试 $B_0=0.3$

系统初始化：

- L1：接入主流AI编程助手
- L2：S(0)=0.8（高支持）F(0)=0.2（低摩擦）
- L3：C(0)=[问题分解:0.3 实现:0.2 调试:0.4 ...]
- L4：策略=新手培养模式

第 1 周：高支持低摩擦

[任务1：实现冒泡排序]

用户请求：「帮我写一个冒泡排序」

L1输出：[完整代码+详细注释]

L2-CFE拦截：

当前F=0.2→应用完整性摩擦

调制后输出：

「这是冒泡排序的思路：

1. 外层循环...
2. 内层循环...

请你先尝试写出框架遇到困难再问我」

[用户尝试]

Alice写出代码但有bug

[L3-AVP-TM] 记录事件

- 任务类型：排序算法
- 完成质量：0.6（有bug但接近正确）
- 独立完成度：0.7（部分独立）
- 更新能力：C\_Alice.implementation += 0.05

[L2-SGS] 削减调整

- 因表现良好准备下周降低S (t)

第 4 周：中支持中摩擦

[当前状态]

- S(t)=0.5（中等支持）
- F(t)=0.5（中等摩擦）
- C\_Alice整体提升至0.5

[任务：实现二叉树中序遍历]

用户请求：「如何写中序遍历」

L1输出：[完整递归代码]

L2-CFE拦截：

应用抽象度摩擦

调制后输出：

「中序遍历的递归思路是：

1. 递归遍历左子树
2. 访问当前节点
3. 递归遍历右子树

基准情况：节点为空时返回

请根据这个思路自己实现代码」

[用户完成任务]

Alice提交代码

[L3-AVP-TM] 评估表现

- 代码质量评分：0.85
- 独立完成度：0.9（几乎无需额外帮助）
- 更新能力：C\_Alice.recursion += 0.1

[L3] 预警更新

- 黄色预警解除（能力提升）
- 记录成功案例

[用户反馈]

系统：

「太棒了您独立完成了中序遍历  
递归能力提升明显当前水平：0.6  
继续保持下个任务难度将略有提升」

### 5.6.2 团队层的 LSA 应用（概念演示）

场景：某软件团队使用 LSA 进行协作开发

L4-MSO 的团队策略：

[团队策略初始化]  
团队规模：8人  
当前T-AVP：0.72（需提升见4.1节）  
目标：T-AVP 0.85（工作假设）

- 策略设计：
- 1. 每周特定时段进入教练模式（L2强制降低S（t））
  - 2. 强制人际代码审查（监测知识流动见4.1节）
  - 3. 新人独立项目测试（L3评估I-AVP）

- [执行]  
特定时段触发：
- 所有成员的AI助手进入"教练模式"  
（只给方向不给代码）
  - 鼓励人际咨询（记录咨询次数）

- [监测]  
L3-AVP-TM每周报告：
- 个体I-AVP通过率：7/8 （87.5%）
  - 团队知识流动指数：0.68（需提升）
  - 识别瓶颈：成员Bob过度依赖AI

- [干预]  
L4触发：
- Bob进入黄色预警
  - 强制参加结对编程（与老手Alice）
  - 下周目标：完成指定数量无AI任务

[结果]  
一段时间后T-AVP复测：0.83（接近目标）

5.7 技术可行性与工程挑战

5.7.1 现有技术栈的适配性

表 5.10: LSA 各层与现有技术的概念映射

LSA 层	核心功能	可用技术方向	成熟度评估
L1	基础 AI 能力	主流大语言模型	高
L2	输出调制	提示工程、微调	中
L2	难度估计	学习分析、IRT 理论	中
L3	能力建模	贝叶斯网络、RL	中-低
L3	行为追踪	遥测系统	高
L4	策略编排	规则引擎、优化	中
全局	多尺度协调	Multi-agent 研究	低

- 关键技术缺口：
- 1. 能力向量的精确建模
- (L3)
- 当前状态：启发式方法、用户自评
  - 需要发展：更精确的认知模型、神经科学启发

## 2. 摩擦强度的自动化调制

(L2)

- 当前状态：手工设计规则
- 需要发展：自适应算法、强化学习方法

## 3. 团队能力的涌现建模

(L4)

- 当前状态：简单聚合方法
- 需要发展：复杂网络理论、多智能体模拟

### 5.7.2 最小可行原型（MVP）设计方向

从完整 **LSA** 到 **MVP** 的简化路径（概念指引）：

核心策略：MVP 强制采用最小遥测 → 逐步丰富策略先观察 AVP 趋势再增加维度（避免过早优化与数据过采）  
从最少的监测指标起步根据实际需求和发现的问题逐步扩展能力向量维度和监测事件类型

第一阶段：单用户基础功能

核心功能：

- L1：集成主流AI模型
- L2-CFE：实现一种摩擦模式（如完整性摩擦）
- L2-SGS：使用预设削减曲线（如指数衰减）

验证目标：

- 用户能完成任务
- 摩擦可接受（通过用户反馈）
- 初步AVP测试通过率达到合理水平

第二阶段：加入监测功能

新增功能：

- L3-AVP-TM：嵌入式能力评估
- 简化能力向量（3 - 5维）
- 基础预警系统（黄色/红色）

验证目标：

- 能检测依赖趋势
- AVP通过率提升
- 预警准确性可接受

第三阶段：团队协调功能

新增功能：

- 团队T-AVP监测
- 简化MSO（团队策略）
- 集体拔线演练支持

验证目标：

- T-AVP测量可行
- 团队策略有效
- O-AVP概念验证

### 5.7.3 关键工程挑战

挑战 1：实时性能

问题：L2/L3 的计算能否在用户可接受延迟内完成

概念解决方向：

- 预计算：提前准备多难度版本
- 异步更新：背景更新  $C(t)$
- 边缘计算：能力向量本地存储

挑战 2：模型对齐

问题：如何让 L1 的 AI 理解适度帮助的语义

概念解决方向：

- RLHF 方法：训练 AI 理解摩擦指令
- 提示工程：设计元提示模板
- 微调方向：在 LSA 场景上进行领域适配

挑战 3：数据冷启动

问题：新用户/新任务如何初始化

概念解决方向：

- 快速评估：简短的能力评估
- 保守初始化：假设较低能力快速调整
- 迁移学习：从相似用户群体借鉴

挑战 4：评估的客观性

问题：如何客观测量用户能力

概念解决方向：

- 多任务交叉验证
- 标准化任务库建设
- 社区验证机制

挑战 5：用户接受度

问题：用户是否愿意接受不完整答案

概念解决方向：

- 渐进引入（初期低摩擦）
- 透明沟通（解释目的）
- 成就系统（能力成长可视化）
- 紧急逃生阀（真正需要时可请求完整帮助）

挑战 6：伦理与隐私

问题：如何在监测能力的同时保护隐私

概念解决方向：

- 联邦学习：本地计算只传聚合
- 差分隐私：添加噪声保护
- 透明审计：用户可查看所有数据

挑战 7：多尺度的复杂性

问题：团队/组织层的协调极其复杂

概念解决方向：

- 分阶段实施（先个体再团队）
- 简化假设（初期忽略复杂交互）



- 人工辅助（L4 初期由管理员操作）

#### 5.7.4 常见问题解答（FAQ）

##### Q1: 权重配置如何确定

A: 权重配置基于场景目标对齐原则需要迭代校准：

- 初始阶段：基于专家判断和领域惯例设定起点
- 校准方法

：

1. 小规模 A/B 测试（参考教育技术领域的实践）
  2. 监测多指标仪表盘
  3. 根据 AVP 通过率、用户留存、能力增长等综合指标调整
  4. 定期复审（如每季度）
- 透明性：所有权重配置及其调整历史应记录支持后续分析

##### Q2: 团队层的 T-AVP 如何量化测量样本量多大才可靠

A: T-AVP 测量的统计考虑（见 4.1 节）：

- 最小样本建议：团队规模 5 人（更小团队建议只做 I-AVP）
- 任务等值性：使用 IRT 校准的平行测验确保测试难度一致
- 置信区间：报告 95% CI 典型需要多次独立测量
- 局限性说明：小样本下估计不稳定建议结合定性观察

##### Q3: MVP 如何开局最小遥测是什么

A: MVP 的最小可行遥测建议：

第一阶段（L1+L2无监测）：

- 记录：任务完成数、用户主动请求帮助次数
- 目的：验证基本可用性

第二阶段（加入L3基础）：

- 新增记录：
  - \* 定期插入无支持任务
  - \* 记录\$P\_2\$分数（可简化为"通过/失败"）
  - \* 用户满意度快速评分
- 目的：初步验证AVP概念

第三阶段（完整遥测）：

- 按表5.6的完整事件集
- 逐步丰富C (t)向量维度

关键：从最少数据开始根据需求渐进增加避免过早优化

##### Q4: 如何处理用户绕过摩擦（使用其他 AI 工具）

A: 这是 LSA 最大的实施挑战之一：

- 技术手段（有限）

：

- 某些情况下可检测其他 AI 工具使用（但易被规避）
- 代码相似度检测（如抄袭检测）
- 文化手段（核心）

:

- 透明沟通 LSA 的目的（能力建构而非限制）
- 成就系统奖励独立完成
- 社区规范（如组织内的” AVP 诚信公约”）
- 设计手段

:

- 让摩擦有意义而非惩罚性
- 提供” 紧急逃生阀”（真正需要时可请求完整帮助）
- 接受现实：完全防止绕过不可能目标是让诚实使用成为主流

**Q5: LSA 是否会让残障用户处于不利地位**

**A:** 这是伦理治理的核心关切（见 5.5.3 节）:

- 原则：“等效努力而非相同任务”（承接 3.0.6 节公平性原则）
- 实践方向

:

- 视障用户：调整任务形式（如口述编程代替键盘输入）但认知摩擦等效
- 学习障碍用户：延长时间、提供辅助工具但 AVP 标准调整为相对提升
- 关键：评估的是能力是否在成长” 而非” 是否达到绝对水平
- 监测：分组 AVP 通过率报告主动识别系统性偏见

**Q6: L1 层可以是闭源模型吗**

**A:** 可以这正是 LSA 分层设计的优势:

- **L1 技术中立**：可以是主流闭源或开源模型
- **L2-L4 独立**：摩擦/监测/编排逻辑不依赖 L1 具体实现
- **可替换性**：随着 AI 技术进步可无缝升级 L1 而不影响上层
- **局限性**：闭源模型的调制能力可能受限于 API 接口
- **开源优势**：如果 L1 是开源模型可以更精细地控制输出

## 5.8 小结：LSA 的理论贡献与实践路径

### 5.8.1 本章核心贡献

#### 1. 架构创新

:

- 首次提出将 CET 理论工程化的完整架构
- 四层分离的设计（L1-L4）使得每层职责清晰支持独立升级和替换
- 明确的层间接口契约（表 5.2）支持多团队并行开发
- 为 AI 辅助工具提供能力建构优先的设计范式

#### 2. 机制设计

:

- **CFE**（认知摩擦引擎）：实现有益摩擦提供多策略空间
- **SGS**（支持削减调度器）：实现渐进独立引入回退与保底机制（承接 3.3.3 节）
- **AVP-TM**（遥测模块）：持续能力监测支持多尺度 AVP（承接 4.5.2 节）
- **MSO**（多尺度编排器）：跨尺度协调整合公平性约束

### 3. 可工程化

:

- 提供从理论到实现的概念路径
- 明确技术栈映射和成熟度评估
- 给出 MVP 实施方向
- 常见问题解答回应实施挑战

#### 5.8.2 与前后章的连接

承接前四章:

- 第三章的 AVP/EML/伙伴式主体性 → LSA 的 L2/L3/L4 实现
- 第三章 3.3.3 节的 S<sub>min</sub> 保底支持 → L2 层的 SGS 回退机制
- 第四章的 T-AVP/O-AVP → LSA 的 L4 多尺度编排
- 第四章 4.5.2 节的多尺度 AVP 监测 → L3 层的 AVP-TM 遥测管线

为第六章铺垫:

- LSA 的每个设计决策都可以证伪
- 明确的技术挑战指向未来研究方向
- 伦理问题需要跨学科讨论

#### 5.8.3 LSA 的三个开放问题

问题 1: 最优摩擦参数是否存在

当前状态: 使用启发式规则 (50-70% 成功率工作假设)

证伪路径:

- 大规模实验: 测试不同  $F(t)$  对 AVP 的影响
- 个体差异: 最优参数是否因人而异
- 领域差异: 编程 vs 写作 vs 数学的最优摩擦是否不同

问题 2: 能力向量能否精确建模

当前状态: 简化的多维向量

挑战:

- 认知能力的维度到底有多少
- 不同能力间的相关性如何
- 能力的动态演化遵循什么规律

未来方向:

- 认知科学 + 机器学习的交叉研究
- 大规模纵向数据收集
- 神经科学的启发 (fMRI EEG)

问题 3: 多尺度协调的理论基础

当前状态: 启发式策略聚合

深层问题:

- 团队能力如何从个体涌现 (复杂系统理论)
- 组织韧性的数学模型 (网络科学)
- 跨尺度的最优控制 (控制理论)

未来方向:

- 多智能体系统研究

- 组织行为学的定量化
- 跨学科的理论整合

#### 5.8.4 实践路径建议

对工具开发者：

1. 从基础功能的 MVP 开始
2. 选择单一领域深耕（如编程、写作）
3. 快速迭代收集真实用户数据
4. 逐步添加监测和编排功能

对研究者：

1. 聚焦某个开放问题（如能力建模）
2. 设计对照实验验证 LSA 假设
3. 发表子领域成果
4. 推动跨学科合作

对组织：

1. 试点项目：选择非关键团队测试
2. 建立 AVP 评估基线
3. 逐步推广到关键业务
4. 建立长期监测体系

本章结束下一章将探讨 CET 理论的局限性、证伪路径与未来研究方向

## 第六章 局限性、证伪路径与未来方向

前五章构建了完整的 CET 理论体系：从理论定位（第一章）、跨学科基础（第二章）、核心机制（第三章）、跨尺度扩展（第四章）到工程化实现（第五章）然而任何有价值的理论都必须明确自身的边界和局限性本章将以学术诚实的态度系统阐述 CET 理论的六大局限、提出清晰的证伪路径、勾勒未来研究议程并讨论伦理与政策考量

本章承诺：

1. 不回避理论的薄弱环节
2. 不过度声称理论的适用范围
3. 提供可操作的证伪路径
4. 主动邀请批判和验证
5. 为后续研究者指明方向

本章逻辑：我们将从理论的内在局限（6.1）出发明确哪些假设可以被证伪（6.2）哪些问题需要未来研究（6.3）最后讨论理论的社会责任（6.4）

### 6.1 理论的六大局限性

#### 6.1.1 尺度边界：个体以下与社会以上

局限陈述：CET 聚焦个体 → 团队 → 组织 → 社会四个尺度但对更微观（神经生理）和更宏观（跨文化/跨代际）的机制涉及不足

具体表现：

1. 神经层面的机制缺失：
  - CET 未深入探讨 AI 使用对大脑可塑性的影响
  - 例如：长期使用 GPS 是否改变海马体？长期依赖 AI 是否改变前额叶？
  - 相关研究（Maguire et al. 2000）显示伦敦出租车司机的海马体变化但 AI 时代的神经科学证据仍然稀缺
  - 后果：CET 的能力退化主张主要基于行为证据缺乏神经机制的直接验证

## 2. 跨文化普适性未验证:

- CET 的案例和论证主要来自 **WEIRD** 社会\* (Western Educated Industrialized Rich Democratic 指西方、受过教育的、工业化的、富裕的、民主的人群与社会语境)
- 集体主义文化 vs 个人主义文化对独立能力的价值判断可能不同
- 例如: 东亚文化中依赖他人可能被视为团队协作美德而非能力缺失
- 后果:AVP 的独立能力标准可能具有文化偏见

## 3. 代际传承的长期机制:

- 第四章讨论了  $T_0/T_1/T_2$  三代模型但跨代传承的具体路径 (家庭/学校/社会) 未充分建模
- 需要 10-20 年的纵向研究才能验证代际鸿沟假说
- 后果:S-AVP 的预测具有高度不确定性

边界声明:

CET的核心适用范围:

- 个体认知能力(可学习、可练习的技能)
- 小团队协作(5 - 50人)
- 单一组织内部(千人以下)
- 10年窗口期内的社会趋势(工作假设)

? 需要谨慎扩展:

- 跨文化应用(需要本地化校准)
- 大规模组织(万人以上复杂度激增)
- 跨代际预测(20年以上不确定性高)

明确不适用:

- 神经生理机制(需要神经科学方法)
- 补偿性外骨骼(如残障辅助目标不同见3.0.6节)
- 纯工具性任务(如计算器不追求能力建构)

### 6.1.2 任务类型的限制

局限陈述:CET 针对认知密集型、可学习的任务但对物理任务、创造性任务、社交任务的适用性有限

适用性矩阵:

表 6.1: CET 在不同任务类型中的适用性

任务类型	适用性	理由	示例
程序性认知任务	高	可明确定义能力、可测量进步	编程、写作、数学
概念性认知任务	中-高	能力可培养但测量难度高	设计、策略规划
创造性任务	中	能力与灵感混合 AVP 难以量化	艺术创作、科学发现
社交任务	低-中	能力受情境影响大难以标准化测量	谈判、领导力
物理任务	低	CET 聚焦认知能力物理技能机制不同	手术操作、运动技能
高风险任务	不适用	拔线测试可能带来不可接受风险	飞行、医疗急救

注: 本表仅用于方向与分层的质量参考不得下推为 *KPI*”; 最终判定仍以 *AVP* 主判据为准” (见 3.0.2 节)

关键限制:

### 1. 能力定义的模糊性:

- 创造性任务中的”能力”难以与”灵感、天赋”分离
- AVP 判据 ( $P_2 - B_0 +$ ) 假设能力可量化但某些能力本质上是定性的

### 2. 测量的生态效度:

- 实验室测量 (如编程测试)vs 真实场景 (如生产环境开发) 的差异
- 高风险、高压环境下的能力表现可能与常规测试大相径庭

### 3. 任务复杂度的上限:

- CET 主要关注中等复杂度任务 (数小时至数天完成)
- 超大规模、跨月的复杂项目 (如大型软件系统设计) 涉及更多组织因素超出 CET 范围

CET 的核心洞察 (认知卸载 → 能力退化) 可能在多种任务中成立但 AVP/EML 的具体操作化需要针对任务类型进行调整本论文提供的测量协议和参数 (如 50-70% 成功率工作假设) 主要基于程序性认知任务的证据

### 6.1.3 测量挑战: 从概念到操作的鸿沟

局限陈述:CET 提出的许多概念 (如  $C(t)$  能力向量、认知摩擦强度  $F$ ) 在理论上清晰但在实践中精确测量极其困难

四大测量难题:

难题 1: 能力向量的维度爆炸

问题:

- 第五章提出  $C(t) = [c_1 \ c_2 \ \dots \ c_n]$
- 但  $n$  应该是多少? 5维? 50维? 500维?
- 不同维度间的相关性如何? 是否存在潜在因子?

当前状态:

- 简化为少数维度 (如 3 - 10 维)
- 但丢失了能力的细粒度信息

理论挑战:

- 认知科学尚未就能力的基本单元达成共识
- 不同理论 (如 Gardner 的多元智能、Cattell 的流体/晶体智能) 给出不同划分

难题 2: 摩擦强度的量化

问题:

- 给出思路但不给代码提供了多少帮助? 60%? 70%?
- 如何客观测量认知努力?

当前近似:

- 信息论: 部分答案的信息熵 / 完整答案的信息熵
- 任务完成预测: 基于历史数据的经验模型
- 用户自评: 主观但快速

局限:

- 信息熵忽略了 "信息有用性" (关键提示 vs 无关细节)
- 经验模型依赖大量数据冷启动困难
- 用户自评存在偏差和不一致性

难题 3: AVP 的等值性保证

问题:

- 如何确保  $T_0$  和  $T$  的任务真正等值?
- 程序性任务 (编程) 相对容易概念性任务 (设计) 极难

当前方法:

- IRT (项目反应理论): 需要大量题库和标定
- 专家评估: 主观且成本高

- 平行测验:需要精心设计

残余风险:

- 任务间的微妙差异可能被忽视
- 用户对任务类型的熟悉度影响表现
- 练习效应与题库污染威胁 (详见表6.2)

难题 4: 长期追踪的流失率

问题:

- AVP需要6 - 12周的纵向追踪 (W=4 - 8周默认6周工作假设需跨领域/任务校准)
- T-AVP/O-AVP需要数月至数年
- S-AVP需要10 - 20年(工作假设)

现实挑战:

- 用户流失率:6周后可能只剩30 - 50%
- 组织变动:关键人员离职、团队重组
- 社会变迁:技术、文化、政策的快速变化

后果:

- 样本偏差(高动机用户过度代表)
- 统计效力下降
- 因果推断困难(混淆变量增多)

表 6.2: AVP 测量的主要效度威胁与缓解策略

威胁类型	具体表现	缓解策略	局限性
等值性威胁	$T_0$ 和 T 任务难度不匹配	IRT 标定、专家双盲评审、平行测验设计	IRT 需大量数据; 专家判断仍有主观性
环境一致性	测试环境与真实使用场景差异	生态效度测试、现场评估、多情境验证	成本高; 难以完全模拟真实压力
练习效应	重复测试导致熟悉度提升	题库轮换、冷却期设计 (2 周)、等值平行卷	题库开发成本高; 冷却期延长研究周期
题库污染	测试题目泄露或过度练习	动态题库、防作弊设计、行为模式检测	无法完全杜绝; 检测算法可能误判
评分者期望	评分者知晓用户使用 AI 而产生偏见	盲评机制、双评分 + 一致性检验、标准化评分 rubric	盲评难度大 (如编程风格暴露); 一致性训练成本高
跨群体公平 (DIF)	不同群体在等值题上出现系统偏差	IRT DIF 检测、项目替换/重标定、分组等值化	需要更大样本与多文化题库; 可能降低可比样本量
延迟巩固威胁	W 期内短期记忆尚未转化为长期能力	W 4 周 (工作假设)、多时间点测量、认知负荷后测	最优 W 因任务而异; 多时间点增加流失率
流失率偏差	高动机用户留存率更高	意向分析 (ITT)、敏感性分析、激励设计	ITT 保守估计可能掩盖真实效应; 激励可能扭曲行为

注: 本表仅用于方向与分层的质量参考不得下推为 KPI; 最终判定仍以 AVP 主判据为准 (见 3.0.2 节)

透明承认:

CET 提供的测量协议 (第三章 AVP、第四章 T-AVP/O-AVP、第五章 LSA 遥测) 代表了当前最佳实践的概念框架而非经过充分验证的标准化工具每个测量协议都需要在特定领域内进行校准、验证、迭代优化我们鼓励研究者将这些协议视为“可证伪的起点而非标准做法”

#### 6.1.4 个体差异与公平性悖论

局限陈述:CET 追求能力提升的普适性目标但个体差异 (学习速度、起点能力、动机) 和公平性要求 (残障、社会经济背景) 之间存在张力

三个悖论:

悖论 1: 标准化 vs 个性化

矛盾:

- AVP需要标准化判据( $P_{20}$   $B_{0}$  + )以确保可比性
- 但个体学习曲线差异巨大(有人3周达标有人3个月)

问题:

- 固定时间窗口 ( $W=4-8$ 周默认6周工作假设) 对快速学习者是浪费对慢速学习者是压力
- 自适应窗口又带来新的公平性问题(如何确定个体最优 $W$ ?)

部分解决:

- 分层标准(新手/中级/专家不同的)
- 自适应窗口+最小/最大约束(4-12周)
- 但本质张力仍存在

悖论 2: 能力 vs 残障

矛盾:

- CET强调独立能力
- 但某些残障使得"独立在传统意义上不可能" (如需要屏幕阅读器)

当前处理(见3.0.6节公平性原则):

- \*\*调整任务形式而不降低挑战强度\*\*
- 评估以\*\*相对提升\*\*而非绝对水平为准
- 核心:"等效努力"原则

未解决的边缘案例:

- 认知障碍 (如学习障碍): 能力提升的速度和上限可能不同
- 神经多样性(如自闭症、ADHD): 最优学习路径可能与典型人群大相径庭
- 临时性障碍(如疾病、压力): 如何处理能力的短期波动?

悖论 3: 机会 vs 结果

矛盾:

- 公平性原则要求机会平等(人人都能通过AVP)
- 但CET设定了通过标准( $P_{20}$   $B_{0}$  + )意味着有人会失败

伦理问题:

- AVP失败者会被贴上依赖标签吗?
- 在组织/教育场景中AVP是否会成为新的能力歧视工具?
- 如何避免AVP强化现有的社会不平等?

设计缓解 (工程化实现):

1. \*\*多次尝试机会\*\*:

- 不公开AVP结果(个人隐私保护)
- 允许多次测试取最佳成绩

2. \*\*提供支持资源\*\*:

- 针对性训练资源
- 能力重建项目(类似再培训)

3. \*\*分组公平性约束\*\* (已在第五章目标函数中体现):

- 优化目标增加硬约束:  
$$\max(\text{AVP通过率}_{\text{群体}}) - \min(\text{AVP通过率}_{\text{群体}}) \leq \text{fairness}$$
- 典型  $\text{fairness} = 0.15$ (15%差异容忍度工作假设)



#### 4. \*\*等效努力的操作化\*\*:

- 残障用户案例: 屏幕阅读器、语音输入
- 守恒的挑战预算: 任务可及性↑但认知挑战度保持不变
- 示例: 视障用户口述编程(等效于键盘输入的认知努力)

#### 5. \*\*核心原则重申\*\*:

计算仅用于标准化与聚合参见第4章跨尺度监测的定位

方法学提醒: 异质性效应与辛普森悖论

在所有涉及分组或亚群的 AVP 评估中必须关注: - 异质性效应 (**HTE**): 报告关键亚组 (如不同年龄、教育背景、起始能力水平) 的 AVP 效应差异 - 辛普森悖论: 检视整体效应与分组效应是否出现反转; 必要时采用分层或多层模型汇总 - 交互作用: 测试干预效应是否因群体特征而显著调节

所有亚组分析必须在预注册中指定避免事后数据挖掘导致的假阳性

### 6.1.5 技术依赖的双刃剑

局限陈述:CET 理论自身依赖 **AI** 技术的稳定性和可访问性这创造了一个潜在的自我指涉悖论

三个依赖层次:

层次 1: 测量依赖

问题:

- AVP/EML的实现依赖AI工具的存在
- 如果AI不稳定 (如API变更、服务中断、成本暴涨)测量本身受影响

风险:

- 厂商锁定: 依赖特定AI供应商(如OpenAI)
- 版本漂移: GPT-4 vs GPT-5的能力差异影响摩擦设计
- 成本波动: 价格变化影响研究可持续性

缓解:

- 供应商多样性(兼容多个API)
- 版本冻结协议(研究期内固定模型版本)
- 开源替代方案(如Llama等)

层次 2: 能力定义依赖

问题:

- 独立能力的定义依赖于AI能做什么
- 当AI能力提升能力边界需要重新划定

示例:

- 2023: 独立编写基础代码被视为重要能力
- 2030 (假设): AI可完全自主编程人类能力转向系统架构?
- 那么AVP测试的对象需要重新定义

理论挑战:

- CET需要认知能力本体论——哪些能力是人类核心价值?
- 但这个本体论本身可能是时代依赖的

层次 3: 社会依赖

问题:

- 大规模AVP评估需要社会资源投入 (教育系统、组织支持)
- 但社会可能认为"与AI协作"比独立能力更重要

潜在冲突：

- 教育政策可能鼓励AI使用而非限制
- 雇主可能更看重使用AI的生产力而非AVP得分
- 代际价值观差异（T\_0代重视独立T\_2代视为过时？）

后果：

- CET可能变成逆时代的理论
- 或者社会最终付出代价时才追悔莫及（代际鸿沟假说见H8）

学术化声明：

此类逆时代类比仅作风险提示与理论边界探讨  
并非价值判断或技术进步的否定CET的核心关切是  
"能力可持续性而非技术使用本身"

#### 6.1.6 文化嵌入性与价值判断

局限陈述:CET 的能力建构目标隐含了特定的价值判断——“独立能力”是值得追求的但这一价值判断并非普适而是文化嵌入的

三个文化维度的挑战：

维度 1: 个人主义 vs 集体主义

西方个人主义文化：

- 重视个体能力、独立性、自我实现
- AVP的独立能力与文化价值契合

东亚集体主义文化：

- 重视团队协作、互相依赖、社会和谐
- 依赖他人不一定被视为缺陷而是团队精神
- AVP可能被视为过度个人主义

可能后果：

- 跨文化实施AVP时遭遇价值观冲突
- T-AVP在集体主义文化中的接受度可能更高（团队层面）

维度 2: 能力观的文化差异

西方分析型认知传统：

- 强调逻辑、规则、独立解决问题
- 适合程序性任务的AVP测试

东亚整体型认知传统：

- 强调情境、关系、整合多方意见
- AVP测试可能无法充分评估这种能力

示例：

- 一个依赖AI做决策支持的东亚管理者
- 西方AVP视角:能力缺失(过度依赖)
- 东亚视角:善用工具、整合资源的智慧

维度 3: 技术关系的哲学差异

西方工具主义传统：

- 技术是达成目的的工具人类是主体
- CET的伙伴式主体性仍保持人的中心地位

东亚和谐观：

- 人与技术共生界限模糊（如茶道中人与茶具的合一）
- 认知外骨骼可能被视为自然演化而非病理

理论后果：

- CET在西方文化中可能更具说服力
- 需要本土化调整以适应不同文化语境

透明承认：

我们承认 CET 理论的价值判断是文化嵌入的主要反映 WEIRD 社会的认知传统这不意味着理论无效但意味着理论的适用边界与普适性声称需要谨慎跨文化验证是未来研究的重要方向（见 6.3.2 节中期研究议程）

## 6.2 八个可证伪假说及其证伪路径

### 6.2.1 核心假说概览与因果识别策略

CET 理论的科学性在于其可证伪性我们明确提出 8 个核心假说并给出清晰的证伪条件这不是为了保护理论免受批评而是为了明确理论的预测边界便于后续研究者检验

元认知原则：我们期待至少部分假说被证伪——这不是失败而是科学进步的标志如果所有假说都被验证那可能意味着我们的预测过于保守或模糊缺乏真正的风险

因果识别策略菜单

CET 假说的验证面临一个核心方法学挑战：如何在无法实施 **RCT** 的情况下建立因果推断？我们提供以下替代策略菜单供研究者根据具体情境选择：

**策略 1：鼓励设计 (Encouragement Design)** - 适用场景：无法强制分配 AI 使用但可以鼓励 - 核心思路：随机分配鼓励使用 AIvs 不鼓励利用工具变量原理 - 示例：向实验组提供免费 AI 订阅对照组不提供 - 分析方法：两阶段最小二乘法 (2SLS) 以“鼓励”作为工具变量 - 优势：保留随机性符合伦理 - 局限：依赖“单调性假设”（鼓励只增加不减少使用）

**策略 2：阶梯楔形设计 (Stepped-Wedge Design)** - 适用场景：组织/团队层面无法同时部署到所有单位 - 核心思路：所有单位最终都接受干预但时间随机错开 - 示例：50 个团队分 5 批引入 EML 工具每批间隔 4 周 - 分析方法：混合效应模型控制时间趋势 - 优势：伦理友好（所有人最终受益）适合实践场景 - 局限：需要假设无时间依赖的干预效应

**策略 3：工具变量 (Instrumental Variables)** - 适用场景：存在影响 AI 使用但不直接影响能力的外生因素 - 核心思路：利用准随机变异进行因果推断 - 示例：公司分配的 AI 许可证数量（受预算而非能力影响） - 分析方法：2SLS 或更鲁棒的估计量 - 优势：可利用观察数据建立因果推断 - 局限：工具有效性难以验证（需要领域知识和敏感性分析）

**策略 4：自然实验 (Natural Experiments)** - 适用场景：外生冲击导致 AI 可用性突变 - 核心思路：利用意外事件（如 API 中断、政策变化）作为准实验 - 示例：\* OpenAI API 大规模宕机事件 \* 某国突然禁止/开放 AI 工具使用 - 分析方法：事件研究法、合成控制法 - 优势：高外部效度真实场景 - 局限：机会稀少难以预先规划可能缺少对照组

**策略 5：回归断点设计 (Regression Discontinuity)** - 适用场景：存在基于阈值的政策/分配规则 - 核心思路：利用阈值附近的准随机分配 - 示例：\* O-AVP 告警阈值 0.70 附近的组织（见 H5）\* 某公司规定绩效 >80 分才能用高级 AI 工具 - 分析方法：局部线性回归检验阈值处的跳跃 - 优势：强因果推断适合政策评估 - 局限：仅对阈值附近群体有效外推需谨慎

选择决策树：

是否能随机分配？

是 → 优先 RCT

否

能否鼓励/激励？→ 鼓励设计

能否错开时间？→ 阶梯楔形

有外生工具? → 工具变量  
有突发事件? → 自然实验  
有阈值规则? → 回归断点

所有准实验方法都依赖不可检验的假设 (如工具有效性、平行趋势) 研究者必须: 1. 明确说明假设及其合理性 2. 进行鲁棒性检验 (如安慰剂测试、敏感性分析) 3. 承认因果推断的局限性 4. 报告所有分析结果 (包括不支持假设的)

表 6.3: 八个核心假说的证据强度与优先级

假说编号	假说名称	证据强度	所需研究设计	预期时间线	优先级
H1	AVP-Basic 假说	Moderate	RCT + 纵向追踪	1-2 年	P0
H2	有益摩擦假说	Preliminary	多臂试验 + A/B 测试	1-3 年	P0
H3	系统性削减假说	Preliminary	对照实验	1-2 年	P0
H4	团队能力极化假说	Strong	自然实验 + 田野研究	2-3 年	P1
H5	组织韧性假说	Moderate	准实验 + 回归断点	3-5 年	P1
H6	摩擦调制假说	Preliminary	A/B 测试 + 用户研究	1-2 年	P2
H7	能力向量假说	Preliminary	降维分析 + 预测建模	2-4 年	P2
H8	代际鸿沟假说	Preliminary	纵向队列研究	10-20 年	P1

注: 本表仅用于方向与分层的质量参考不得下推为 *KPI*; 证据强度和时间线均为工作假设需根据实际研究进展校准

证据强度说明: - **Strong**: 多个独立研究收敛至少 1 个高质量 RCT 或准实验 - **Moderate**: 部分实证支持但样本量小或设计局限 - **Preliminary**: 主要基于理论推导和方向性观察证据

表 6.4: 核心假说的证伪对照与理论修正路径

假说	最小样本/时长	核心证伪观测	因果识别策略	失败后的修正方向
H1	N 200 W=4-8 周 (默认 6 周)	$P_2-B_0$ 效应量 $<0.2$ 或反向	RCT 或鼓励设计	检查 EML 条件是否充分执行; 重新校准 阈值; 可能放弃 AVP 作为唯一标准
H2	N 300 12 周	无倒 U 曲线或峰值在 $<30\%/>85\%$	多臂试验	承认” 最优区间因任务而异”; 改为自适应个性化策略; 可能放弃统一摩擦参数
H3	N 240 12 周	固定支持 削减组	对照实验	承认削减可能非必要; 改为按需支持” 模型”; 保留 AVP 但修正 EML
H4	50-100 团队 6-12 月	能力方差无变化或缩小	自然实验 + DID	承认团队协作可能缓冲极化; T-AVP 可能更依赖组织文化而非技术设计
H5	20+ 组织 12-24 月	O-AVP 与恢复时间无关	回归断点 + 事件研究	承认 O-AVP 双阈值模型可能不普适; 改为连续风险评估; 可能放弃单一阈值模型
H6	N 300 12 周	动态摩擦无优势或劣于固定	A/B 测试	承认简单即美; 删除自适应复杂度; LSA 简化为固定策略库
H7	N 1000 6-12 月	能力向量解释方差 $<10\%$	降维 + 预测建模	承认能力可能高维不可约; 改为任务特定能力模型; 放弃统一 $C(t)$ 向量
H8	3 代队列 15-20 年	Cohens $d<0.3(T_2$ vs $T_0)$	纵向队列	承认代际差异可能被高估; 教育干预有效; CET 警示过度但原则仍有价值

注: 本表仅用于方向与分层的质量参考不得下推为 *KPI*; 修正方向为理论演化预案并非失败即放弃理论

修正方向说明: - 参数校准型: H2、H5 - 核心原则保留调整具体参数 - 条件精细化型: H3、H4 - 识别假说成立的边界条件 - 模型简化型: H6 - 承认复杂度不必要回归简单设计 - 理论修正型: H1、H7 - 可能需要修改核心假设 - 警示价值型: H8 - 即使被证伪警示作用仍有意义

### 6.2.2 个体层假说的证伪路径

#### 假说 H1: AVP-Basic 假说

声明:

在程序性认知任务中通过有益摩擦 ( $F=50-70\%$ 成功率工作假设)

+ 系统性支持削减( $S4 \rightarrow S1 \rightarrow S0$ )设计的AI工具协作W周后

用户在拔线窗口内的独立表现将显著优于基线( $\$P\_2\$ - \$B\_0\$ + 0.3 SD$ )

证伪条件:

- 在严格RCT设计下对照组(传统学习)与实验组(EML设计)在 $\$P\_2\$$ 表现上无显著差异(效应量  $< 0.2$ )
- 或者实验组 $\$P\_2\$ < \$B\_0\$$ (能力退化)

验证方法:

- 设计:  $2 \times 2$ 因子RCT(有益摩擦 $\times$ 支持削减)
- 样本: N 200(每组50)多领域复制(编程/写作/数学)
- 测量: 标准化AVP协议(见第三章)双盲评分
- 分析: 意向分析(ITT) + 符合方案分析(PP)
- 控制: 等值性验证、流失率分析、敏感性分析

执行忠实度: 记录EML条件的达成率(如削减曲线按计划执行比例、观察到的实际成功率是否落在50-70%区间)

操纵检查: 对关键组别进行独立抽样复核确认摩擦强度与支持档位与实验指定一致(避免名义分组、实则偏离)

数据收集:

- $\$T\_0\$$ :  $\$B\_0\$$ 基线测试
- $\$T\_1\$$ - $\$T\_2\$$ : 协作期 ( $W=4-8$ 周默认6周可校准)
- $\$T\_3\$$ :  $\$P\_2\$$ 拔线测试 ( $W=4-8$ 周默认6周后工作假设需跨领域/任务校准)
- $\$T\_4\$$ : 长期保持测试 (6个月后可选)

预期结果:

- 实验组(EML):  $\$P\_2\$ - \$B\_0\$ = 0.3 SD$  (工作假设)
- 对照组(无摩擦/无削减):  $\$P\_2\$ - \$B\_0\$ = 0$  或  $< 0$

如果Cohens  $d < 0.2$ 或方向相反则H1被证伪

短期试点建议:

可先进行24h拔线轻量测(I-AVP级)合格后再进入48h/长窗;

流程与评分表可参考附录轻量SOP

#### 假说 H2: 有益摩擦假说

声明:

存在一个"最优挑战区" (50-70%成功率工作假设)在此区间内

用户的能力提升( $\$P\_2\$ - \$B\_0\$$ )最大化过高摩擦(成功率 $<30\%$ )

导致挫败过低摩擦(成功率 $>85\%$ )导致卸载

证伪条件:

- 证明摩擦强度与能力提升呈线性关系(无倒U曲线)
- 或证明最优区间显著偏离50-70%(如最优点在20%或90%)

验证方法:

- 设计: 多臂试验5-7个摩擦水平(20%/40%/60%/80%/95%)
- 样本: N 300(每臂50)跨任务类型复制
- 测量: 连续变量( $\$P\_2\$ - \$B\_0\$$ )用户满意度坚持率

- 分析:二次回归、分段回归、贝叶斯模型比较

预期模式:

- 倒U型曲线峰值在50 - 70%区间(工作假设)
- 如果线性递增或峰值在<30%或>85%则H2被证伪

校准参数:

- 不同任务类型可能有不同的最优区间
- 需要领域特定的校准和敏感性分析

假说 **H3**: 系统性削减假说

声明:

系统性支持削减(S4→S1→S0→0)优于固定支持渐进削减组的  
AVP通过率和长期保持显著高于固定支持组

证伪条件:

- 固定支持组(如持续S3)的\$P\_2\$表现不劣于削减组
- 或者削减速度与结果无关(线性/指数/S型无显著差异)

验证方法:

- 设计:3×2因子实验(削减模式×削减速度)
  - \* 削减模式:线性、指数、S型
  - \* 削减速度:快(4周)、慢(12周)
  - \* 对照:固定S3支持
- 样本:N 240(每组40)
- 测量:AVP通过率、长期保持(6个月)、用户体验

执行忠实度:记录每周实际支持档位确认削减曲线执行与设计一致;记录用户感知的支持变化轨迹

预期结果:

- 削减组>固定组 (Cohens d 0.3 或 10% (\*working assumption\*))
- S型削减可能最优(但需验证)

如果固定支持组同等或更好则H3被证伪

### 6.2.3 团队与组织层假说的证伪路径

假说 **H4**: 团队能力极化假说

声明:

在无EML约束的AI使用下团队内部出现能力极化:  
高能力者进一步提升低能力者能力退化T-AVP下降

证伪条件:

- 证明团队内部能力方差无显著变化
- 或证明低能力者也获得提升(无极化效应)

验证方法:

- 设计:自然实验(AI引入前后对比)
- 样本:50 - 100个软件开发团队
- 测量:个体I-AVP、团队T-AVP、知识流动网络
- 分析:方差变化、基尼系数、社会网络分析

数据收集:

- \$T\_0\$: 团队引入AI前的基线

- $T_1-T_2$ : 6 - 12个月使用期
- $T_3$ : 周五无AI日拔线测试

预期模式:

- 能力方差显著增大 (F检验  $p < 0.05$ )
- 基尼系数增加 0.15(工作假设)
- 知识依赖网络出现中心化

如果无极化或负向极化则H4被证伪

假说 **H5**: 组织韧性假说 (双阈值模型)

声明(修订版):

O-AVP采用双阈值模型:

- 告警阈值: 0.70(工作假设触发风险排查/回退)
- 目标阈值: 0.85(工作假设用于质量分层/优秀实践识别)

双阈值的定义与计算口径见4.2.2节(O-AVP公式)本章仅做阈值使用与证伪路径说明

$O-AVP < 0.70$  (告警阈值) 的组织在AI中断后恢复时间显著更长  
面临显著的系统性风险

证伪条件:

- 找到  $O-AVP < 0.70$  (告警阈值) 但在AI中断后快速恢复(<12h)的组织
- 或证明阈值不具有统计上的显著性

验证方法(伦理约束下):

- 方案A: 自愿参与的48h演练(第四章协议)
- 方案B: 自然实验(AI服务商宕机事件的事后分析)
- 方案C: 仿真模拟(Agent-Based Model)

数据收集:

- 记录20+组织的O-AVP得分
- 追踪真实AI中断事件的恢复时间
- 分析: O-AVP与恢复时间的关系
- 预测:  $O-AVP < 0.70$  (告警阈值) 的组织恢复时间显著更长

备选因果策略:

- 回归断点设计: 利用0.70阈值前后的跳跃
- 差分中的差分: AI宕机前后的组织对比

统计考虑(团队/组织层):

- 按集群/分层设计估算样本量(考虑ICC)
- 使用多层线性模型或GEE做稳健推断
- $\alpha=0.05$ 、Power 0.8(工作假设)需敏感性分析

口径声明(工作假设):

O-AVP双阈值模型:

- 告警阈值: 0.70 (触发风险排查/回退)
- 目标阈值: 0.85 (质量分层/优秀实践识别)
- (工作假设需跨领域/任务校准)

不同领域需做敏感性分析与校准

#### 6.2.4 LSA 设计假说的证伪路径

假说 **H6**: 摩擦调制的有效性

声明:

L2层的动态摩擦校准引擎(CFE)优于固定摩擦

证伪条件:

- 证明固定摩擦的效果不差于动态调整
- 或证明动态调整的成本超过收益

实验设计(A/B测试):

- 对照组: 固定 $F=0.6$ (中等摩擦)
- 实验组: 动态CFE(根据成功率调整)
- 测量: 12周后的AVP通过率、用户满意度、系统成本
- 预测: 实验组AVP通过率更高且满意度不显著降低

执行忠实度: 记录CFE的实际调整轨迹确认是否按成功率反馈动态调整

如果对照组同等或更好则H6被证伪

理论意义:

- 如果被证伪意味着简单即美
- 复杂的自适应系统可能是过度设计

假说 **H7**: 能力向量的可建模性

声明:

用户的认知能力可以用低维向量 (<20维)有效表征

证伪条件:

- 证明能力本质上是高维、非线性、不可压缩的
- 或证明能力向量无法预测独立表现

实验设计:

- 收集1000+用户的行为数据
- 用降维算法(PCA VAE)提取能力向量 $C(t)$
- 测试:  $C(t)$ 能否预测用户在新任务上的表现
- 预测: 解释方差>30%

如果<10%则H7被证伪

认知科学背景:

- 心理测量学的因子分析传统
- 但AI时代的能力结构可能与传统不同

#### 6.2.5 社会层假说的证伪路径

假说 **H8**: 代际能力鸿沟

声明:

\$T\_2\$代(2015年后出生AI原生代)在无AI下的独立能力  
将显著低于\$T\_0\$代(1980 - 2000年)

证伪条件:

- 2035 - 2040年的纵向数据显示:  
\$T\_2\$代与\$T\_0\$代的能力差异不显著(Cohens  $d < 0.3$ )



实验设计(纵向队列研究):

- 从2025年开始追踪:
  - \*  $T_0$ 代(当前25 - 45岁): 基线能力测试
  - \*  $T_1$ 代(当前10 - 25岁): 中期测试
  - \*  $T_2$ 代(当前0 - 10岁): 持续追踪到2040年
- 测量: 标准化能力测试(如数学、阅读、问题解决)
- 控制变量: 教育水平、社会经济地位
- 预测:  $T_2$ 代在无AI条件下的表现 <  $T_0$ 代

如果差异不显著则H8被证伪

时间尺度:

- 需要15 - 20年才能获得决定性证据
- 是CET理论中最长期的预测

S-AVP定位澄清:

S-AVP产出的是人群层长期趋势的可证伪预测  
不直接等同于拔线实验的判定;用于指导社会级  
"研究议程/政策试探"非组织/个体的即时评分

元认知反思:

上述 8 个假说代表 CET 理论的核心可证伪预测它们不是全部但是最关键的我们热切期待其中一些假说被证伪——理论的进步往往来自于发现我们错在哪里如果所有假说都被验证那可能意味着我们的预测过于保守或模糊真正有价值的理论应该大胆猜测并提供清晰的证伪路径

## 6.3 未来研究议程: 三个时间尺度

### 6.3.1 短期研究 (1-3 年): 验证核心机制

优先级 1: AVP 协议的标准化

目标: 将第三章的 AVP 协议转化为可复现的测量工具

关键任务:

1. 跨领域校准:
  - 在编程、写作、数学、语言学习等 5 个领域实施 AVP
  - 确定各领域的  $\alpha$ 、 $\beta$ 、 $\gamma$ 、 $\delta$ 、 $\epsilon$ 、 $\zeta$ 、 $\eta$ 、 $\theta$ 、 $\iota$ 、 $\kappa$ 、 $\lambda$ 、 $\mu$ 、 $\nu$ 、 $\xi$ 、 $\omicron$ 、 $\pi$ 、 $\rho$ 、 $\sigma$ 、 $\tau$ 、 $\upsilon$ 、 $\phi$ 、 $\chi$ 、 $\psi$ 、 $\omega$  摩擦参数
  - 发表《AVP 测量手册》(类似 IRT 测量标准)
2. 信度与效度验证:
  - 重测信度: 同一用户多次测量的一致性
  - 效标效度: AVP 是否预测真实场景表现
  - 构念效度: AVP 是否真正测量独立能力
3. 开源工具包:
  - 发布 AVP 测量软件 (支持自动化测试)
  - 提供任务库 (至少 100 个标准化任务)
  - 建立社区数据共享平台

预期产出: - 3-5 篇实证论文 (各领域的 AVP 验证) - 开源 Github 项目 (星标 >1000) - 至少 5 个独立团队复现

预注册与开放科学: - 重要实验预注册 (OSF/AsPredicted) - 负结果也入库 (避免发表偏倚) - 提供最小可复现包 (题本、评分 Rubric、统计脚本)

全球验证机会: 全球范围的实证研究为 CET 理论提供了重要的验证机会。对五国教师和学生的调查显示缺乏适当引导的 AI 使用往往导致学习者追求效率而牺牲深度理解 (Aguilar et al. 2025)。这强调了 EML 设计中渐进撤出原则的重要性, 也为 H3 假说的跨文化验证提供了初步证据。未来研究应在不同文化和教育体系中系统测试 CET 的普适性。

#### 优先级 2: EML 参数的实验优化

目标: 确定 50-70% 成功率、“S4→S1→S0 削减是否最优” (工作假设)

关键任务:

1. 摩擦参数实验:
  - 5×5×5 因子设计: 成功率区间、削减速度、任务类型
  - N=500-1000 用户 12 周追踪
  - 测量: AVP 通过率、学习曲线、用户满意度
2. 削减曲线优化:
  - 对比: 线性、指数、S 型、阶梯式削减
  - 识别过快削减的预警信号
  - 个性化削减策略的收益分析
3. 多模态摩擦:
  - 探索: 完整性摩擦 vs 抽象性摩擦 vs 延迟性摩擦
  - 不同任务类型的最优摩擦组合

预期产出: - 数据驱动的 EML 参数指南 - 摩擦设计模式库 - 可能推翻或修正当前参数

#### 优先级 3: 小规模 LSA 原型

目标: 实现 L1+L2+L3 的 MVP 在单一领域验证

关键任务:

1. 选择垂直领域:
  - 建议: 编程教育 (数据易获取、能力易测量)
  - 备选: 写作辅助、数学学习
2. 开发最小功能:
  - L1: 集成 GPT-4 API
  - L2: 固定摩擦模式 (完整性摩擦)+ 预设削减曲线
  - L3: 嵌入式微测试 + 黄色/红色预警
3. 用户研究:
  - 招募 100-200 名用户
  - 6 个月纵向追踪
  - 收集定量 (AVP 通过率)+ 定性 (访谈) 数据

预期产出: - 工作原型 (可演示) - 案例研究论文 - 识别实施中的关键挑战

### 6.3.2 中期研究 (3-5 年): 跨尺度扩展与理论整合

方向 1: 团队与组织层的实证研究

目标: 验证第四章的 T-AVP/O-AVP 框架

关键任务:

1. **T-AVP** 大规模研究:
  - 与 50-100 个软件团队合作

- 实施”周五无 AI 日” + 集体拔线演练
- 追踪团队知识流动、协作模式、T-AVP 得分
- 识别团队层面的成功模式

## 2. O-AVP 中断模拟:

- 与 20-50 个组织合作
- 实施 48 小时 AI 中断演练 (伦理委员会批准)
- 测量恢复时间、应急能力、组织韧性
- 验证 O-AVP 与恢复速度的关系

## 3. 跨尺度机制研究:

- 追踪 I-AVP→T-AVP→O-AVP 的传导路径
- 识别”涌现性断裂”(个体通过但团队失败)
- 建立多层次因果模型

预期产出: - T-AVP/O-AVP 的标准化协议 - 组织能力建构的最佳实践手册 - 跨尺度理论的精细化模型

## 方向 2: 跨文化适应性研究

目标: 验证 CET 在非 WEIRD 文化中的适用性

关键任务:

### 1. 概念等值性验证:

- 独立能力在不同文化中的语义差异
- AVP 判据的文化公平性评估
- 需要哪些本土化调整?

### 2. 多文化对照研究:

- 至少 3 个文化群体 (如北美、东亚、拉美)
- 相同任务、相同协议
- 比较 AVP 通过率、用户体验、长期效果

### 3. 文化特定机制探索:

- 集体主义文化下的 T-AVP 是否更有效?
- 高情境文化下的有益摩擦定义是否不同?
- 识别普适性 vs 文化特异性的边界

预期产出: - 跨文化 AVP 校准指南 - 文化嵌入性理论的精细化 - 可能推翻或修正普适性声称

## 方向 3: 神经科学整合

目标: 为 CET 提供神经层面的验证

关键任务:

### 1. AI 使用的神经影响:

- fMRI 研究: 长期 AI 使用 vs 非使用者的大脑差异
- 是否有类似”GPS 效应”(海马体萎缩) 的证据?
- 有益摩擦 vs 零摩擦的神经激活模式差异

### 2. 能力退化的神经标记:

- 识别能力退化的早期神经信号
- 建立神经 → 行为的预测模型
- 为 AVP 测试提供生物标记补充

### 3. 可塑性窗口研究:

- 不同年龄段的能力建构可塑性

- $T_0/T_1/T_2$  代的神经差异
- 为代际鸿沟假说提供神经证据

预期产出: - CET 的神经科学基础 - 能力退化的生物标记 - 跨学科整合理论

### 6.3.3 长期研究 (5-10 年以上): 社会影响与理论演化

方向 1: 代际纵向研究

目标: 验证 S-AVP 和代际鸿沟假说 (H8)

关键任务:

1.  $T_0/T_1/T_2$  队列建立:
  - 从 2025 年开始追踪 3 个代际队列
  - 每 2 年一次标准化能力测试
  - 控制教育、社会经济地位等混淆变量
2. 关键时间点测量:
  - 2030:  $T_1$  代 (15-30 岁) 中期评估
  - 2035:  $T_2$  代 (10-20 岁) 初步评估
  - 2040: 三代对比验证代际鸿沟假说
3. 社会层面干预研究:
  - 政策实验: 引入 EML 原则的教育改革
  - 对照设计: 实验区 vs 对照区
  - 评估: 干预是否缓解代际鸿沟

预期产出: - 15-20 年的纵向数据集 - 代际能力演化的决定性证据 - 为教育政策提供实证基础

方向 2: AI 能力演化的理论适应

目标: 随着 AI 能力提升调整 CET 理论边界

关键任务:

1. 能力本体论的动态更新:
  - 每 5 年重新定义核心人类能力
  - 区分永恒价值能力 vs 时代依赖能力
  - 建立能力分类的动态框架
2. AGI 情境的理论推演:
  - 如果出现通用人工智能 CET 如何适应?
  - “独立能力” vs “增强能力”的边界重新讨论
  - 可能需要理论范式转移
3. 理论演化的三层同心圆:
  - 内层 (核心原则): AVP/EML 的适用性验证
  - 中层 (参数):  $\alpha$ 、 $W$ 、摩擦参数的持续校准
  - 外层 (实现): LSA 的技术更新

预期产出: - CET 2.0 理论 (如果需要根本性修正) - 动态能力本体论框架 - 理论演化的元模型

方向 3: 开放议题探索

议题 A: S 层的大过滤器脆弱性 (见第四章 4.4.3 节)

问题:

- 是否存在文明层面的认知退化风险?
- 代际鸿沟是否构成长期生存脆弱性?

研究方向：

- 跨文明比较研究(如果有足够样本)
- 历史案例分析(技术依赖导致的文明脆弱性)
- 复杂系统建模(认知公地悲剧的临界点)

透明声明：

这是高度推测性的长期风险CET不做确定性结论  
仅将其列为研究议程以供探索作为可证伪迹象

议题 **B**: 人机融合的哲学边界

问题：

- 当脑机接口出现"独立vs"增强的边界何在？
- 是否存在纯粹人类认知vs混合认知的本质区别？

研究方向：

- 哲学论证：主体性的边界
- 实证研究：脑机接口用户的能力演化
- 伦理讨论："认知增强"的道德限度

理论意义：

可能需要CET的根本性修正或范式转移

## 6.4 研究伦理与开放科学承诺

### 6.4.1 伦理原则

CET 理论的验证涉及人类受试者必须严格遵守研究伦理：

1. 知情同意：所有 AVP 测试、拔线演练必须获得参与者知情同意
2. 无伤害原则：拔线测试不得用于高风险任务 (如医疗、飞行)
3. 隐私保护：AVP 结果个人隐私不得用于雇佣/教育歧视
4. 公平性原则：
  - 为残障、弱势群体调整任务形式而不降低挑战强度
  - 评估以相对提升而非绝对水平为准
  - 守恒的挑战预算：任务可及性  $\uparrow$  但认知挑战度保持不变
5. 撤回权：参与者可随时退出研究不受惩罚

### 6.4.2 开放科学承诺

为促进理论验证和批判我们承诺：

1. 数据开放：
  - 匿名化数据集公开发布 (符合隐私法规)
  - 原始数据存储于开放平台 (如 OSF)
  - 对敏感场景可采用差分隐私或联邦学习的开源实现以降低再识别风险
2. 方法透明：
  - 详细研究协议预注册
  - 在预注册中指定主要终点 (如  $P_2-B_0$  效应量) 与次要终点 (保持/迁移/成本) 并提交多重比较校正方案避免事后指标淘选
  - 统计代码开源 (GitHub)
  - 负结果报告 (避免发表偏倚)
3. 工具开源：

- AVP 测量软件开源
- LSA 参考实现开源
- 题库与评分 rubric 公开

#### 4. 协作邀请:

- 欢迎独立团队复现
- 鼓励跨文化验证
- 接受批评性检验

## 6.5 结语：理论的生命在于批判与演化

CET 理论诞生于 2025 年——AI 能力爆发、人类认知面临重构的关键时刻我们提出这个理论不是因为我们相信它是“完美的或最终的”而是因为现在迫切需要一个可证伪的、系统性的框架来理解和引导人机共生的未来

本章揭示的六大局限提醒我们：CET 是在特定技术、文化、认识论背景下的产物它的价值不在于永恒正确而在于：

1. 提供可证伪的预测：8 个核心假说都有明确的证伪条件
2. 承认不确定性：所有参数都标注为工作假设需校准
3. 邀请批判：我们期待被证伪而非害怕被证伪
4. 指明研究方向：三个时间尺度的研究议程为后续工作者铺路
5. 保持演化能力：三层同心圆架构允许理论随证据更新

最后的呼吁：

如果你是研究者请：- 挑战 CET 的假说用严格的实证研究证伪或验证 - 在不同文化、不同领域复制 CET 的核心发现 - 提出竞争性理论推动领域进步

如果你是开发者请：- 将 EML 原则融入 AI 工具设计 - 测量并公开你的产品的 AVP 表现 - 参与开源社区共建能力建构优先的 AI 生态

如果你是教育者/管理者请：- 在组织中试点 AVP 评估 - 关注团队和组织的认知韧性 - 平衡效率与能力建构的长期价值

如果你是政策制定者请：- 关注 CET 揭示的长期风险（代际鸿沟、认知公地悲剧）- 支持跨学科的纵向研究 - 建立 AI 工具的能力建构影响评估机制

科学理论不是圣经而是工具 CET 的最大价值不在于给出答案而在于“提出正确的问题”我们相信即使 CET 的某些假说最终被证伪它也已经完成了使命——推动我们更深入地思考人类与 AI 共存的未来

理论的生命在于被讨论、被检验、被超越 我们期待那一天的到来

## 第 7 章 - 术语与符号系统

文档版本：v1.0 | 已应用补丁包 v1.0 用途：核心概念速查、参数一致性维护、术语标准化 性质：只读参考文档（唯一修订源在主文相应章节）

### 1. 核心概念固定锚点（B1-B5）

使用说明：以下内容第 3.0.2-3.0.6 节之逐字转录唯一修订源在第 3 章如需更新请先修改 3.0 节本处为只读副本供快速查阅使用

#### B1 | AVP 定义锚点

反脆弱性验证原则（Antifragility Validation Principle AVP）：以拔线测试检验协作是否促进独立能力

判据： $P_2 \geq B_0 +$

其中：-  $B_0$ ：使用 AI 前的独立基线能力 -  $P_2$ ：协作一段时间后在拔线窗口  $W=4-8$  周（默认 6 周）内的独立表现 -  $Cohens\ d\ 0.3$  或  $10\%$ （*working assumption*）-  $P_1$ （协作期表现）不参与最终判定 \*\*

**Canonical source:** 见 3.0.2 节（本处为只读副本）

**B2 | EML 定义锚点**

内共生最小法则（**Endosymbiotic Minimal Law EML**）：构成认知内共生的设计必要条件为：

- (1) 有益认知摩擦：使用户处于最优挑战区（群体级工作假设成功率 50-70% 需跨领域/任务校准个体自适应）
- (2) 系统性支持削减：AI 支持强度按既定削减曲线从  $S4 \rightarrow S1 \rightarrow S0$

二者为联合充分的设计条件但最终仍需 **AVP** ( $P_2 - B_0 +$  ) 作为验收必要条件

**Canonical source:** 见 3.0.3 节（本处为只读副本）

**B3 | LSA-F 功能分层锚点**

**LSA-F**（功能分层）：- **L1** 知识整合层：基础 AI 能力接入 - **L2** 状态建模层：摩擦设计与支持削减 - **L3** 摩擦校准层：能力监测与预警 - **L4** 元认知协调层：多尺度编排与伦理治理

支持档位栈（ $S4 \rightarrow S1 \rightarrow S0$ ）用于表达支持强度与 LSA-F 为正交维度

**Canonical source:** 见 3.0.4 节（本处为只读副本）

**B4 | 最优挑战区锚点**

最优挑战区：为促成长期保持与迁移系统应将任务难度/提示强度自适应调至成功率 **50-70%**（工作假设需跨领域/任务校准随任务与个体校准）>85% 近似卸载、<30% 易致挫败

**Canonical source:** 见 3.0.5 节（本处为只读副本）

**B5 | 边界条件锚点**

边界条件：本理论适用于能力增强型人机协作补偿性外骨骼（如残障辅助、超越生理极限的设备）不适用此判据所有参数均为概念工作模型需跨领域校准

**Canonical source:** 见 3.0.6 节（本处为只读副本）

**2. 核心参数速查表**

参数名称	符号/公式	默认值/范围	理论依据	校准方向
最小提升 阈值		Cohens d 0.3 或 10% ( <i>working assumption</i> )（工作假设需跨领域/任务校准）	Cohens d 中等 效应量	程序性任务可能更低（0.2-0.3）创造性任务可能需更高（0.4-0.5）
拔线窗口	W	4-8 周（默认 6 周）（工作假设需跨领域/任务校准）	能力巩固的经验性时间窗口	快速技能 4 周复杂技能 8-12 周长期学习 12-24 周
最优挑战区	成功率区 间	50-70%（工作假设需跨领域/任务校准）（工作假设需跨领域/任务校准）	心流理论 + 最近发展区 (ZPD)	个体自适应任务特定校准文化敏感性调整
支持削减 起点	S0	0.8（80% 支持）（工作假设需跨领域/任务校准）	平衡初期依赖与独立能力培养	高风险任务可从 0.6 开始新手友好任务可从 0.9 开始

参数名称	符号/公式	默认值/范围	理论依据	校准方向
安全支持下限	S_min	0.2（不低于 20%）（工作假设需跨领域/任务校准）	防止用户完全迷失的安全网	专家用户可降至 0.1 新手建议不低于 0.3
削减速率		按任务调整（工作假设需跨领域/任务校准）	个体学习曲线差异	线性/指数/S 型曲线择优
团队 AVP 阈值	T-AVP	0.7（群体级）（工作假设需跨领域/任务校准）	组织韧性研究	不同团队规模需调整关键任务建议 0.8
组织 AVP 阈值	O-AVP	告警 0.70 目标 0.85（工作假设需跨领域/任务校准）	中断演练恢复能力	按行业风险容忍度调整 48h 窗口可改为 24h/72h
代际窗口期	$T_0 \rightarrow T_1$	10 年（2025-2035）（概念占位符）	技术代际影响推测	纵向研究校准跨文化验证

注：所有参数均为概念工作模型需通过实证研究跨领域校准不同任务类型、用户群体、文化背景可能需要显著调整

3. 缩写速查表

缩写	全称	中文	首次出现
AVP	Antifragility Validation Principle	反脆弱性验证原则	1.3 节
EML	Endosymbiotic Minimal Law	内共生最小法则	1.3 节
LSA	Layered Symbiosis Architecture	分层共生架构	1.3 节
LSA-F	LSA Functional Hierarchy	LSA 功能分层	3.0.5 节
CFE	Cognitive Friction Engine	认知摩擦引擎	5.2 节
SGS	Support Graduation Scheduler	支持削减调度器	5.3 节
AVP-TM	AVP Telemetry Module	AVP 遥测模块	5.4 节
MSO	Multi-Scale Orchestrator	多尺度编排器	5.5 节
I-AVP	Individual AVP	个体反脆弱性验证	4.1 节
T-AVP	Team AVP	团队反脆弱性验证	4.1 节
O-AVP	Organizational AVP	组织反脆弱性验证	4.2 节
S-AVP	Societal-level indicators	社会层认知资本指标	4.3 节
BCI	Business Continuity Index	业务连续性指数	4.2 节
ICR	Independent Completion Rate	独立完成率	4.2 节
IRT	Item Response Theory	项目反应理论	附录 A
RCT	Randomized Controlled Trial	随机对照试验	6.2 节
DID	Difference-in-Differences	差分中的差分	6.2 节
ZPD	Zone of Proximal Development	最近发展区	2.3 节

4. 核心评估原则

4.1 等效努力原则（Equivalent Effort Principle）

在 AVP 测试中对不同能力水平或特殊需求的用户应遵循等效努力”而非”等值任务原则：

三句口径（唯一标准表述）：

- 1. 调整任务形式而不降低挑战强度
  - 示例：视力障碍者可使用语音测试但问题难度保持不变
  - 原则：改变呈现方式不改变认知负荷



- 2. 评估以相对提升为准
  - 判据： $P_2 \geq B_0 + \Delta$ （相对于个人基线的提升）
  - 而非： $P_2$  某个绝对标准
  - 理由：尊重个体差异避免一刀切

- 3. 挑战预算守恒
  - 核心：认知负荷总量保持一致
  - 方法：通过任务分解、时间调整、辅助工具等方式平衡
  - 目标：确保不同用户面临等效的认知挑战

应用场景：- 残障人士的能力评估（见 3.1.4 节）- 跨文化任务等值性校准（见 6.1.1 节）- 不同年龄群体的适应性调整 - 教育背景差异的补偿设计

引用说明：其他章节提及等效努力时使用（见第 7 章 § 4.1 等效努力原则）引用此处

#### 4.2 Goodhart 防护原则

核心问题：当度量成为目标时它就不再是好的度量（Goodharts Law）

**CET** 的应对策略：

- 1. **AVP** 判据的非 **KPI** 化
  - AVP 分级（Basic/Retention/Transfer）仅用于质量分层
  - 禁止将 AVP 分数用于人事考核、绩效排名、资源分配
  - 任何涉及利益分配的场景都不得使用 AVP 作为唯一判据
- 2. 固定脚注模板（在所有阈值表格下使用）：> 注（*Goodhart* 防护）：本表/分级仅用于方向与分层不得下推为 **KPI** 最终判定以 **AVP** 主判据（见 3.0.2 节）为准
- 3. 监测与预警的分离
  - 监测数据（如  $C(t)$  能力向量）仅用于系统改进
  - 不与个人利益挂钩
  - 匿名化处理保护用户隐私

反面案例（见 CET9 附录 A FAQ Q5）：某公司将 AVP 用于晋升评估导致用户人为操纵基线、拔线期违规使用 AI 完全失效

### 5. 核心术语中英对照

#### 5.1 理论核心概念

中文术语	英文术语	缩写	核心定义（简版）	详见
认知内共生	Cognitive Endosymbiosis	-	AI 作为伙伴通过摩擦与削减促进能力提升	3.0.4 节
认知外骨骼	Cognitive Exoskeleton	-	过度依赖 AI 导致独立能力退化的病理模式	3.0.6 节
反脆弱性验证原则	Antifragility Validation Principle	AVP	通过拔线测试验证能力是否提升	3.0.2 节
内共生最小法则	Endosymbiotic Minimal Law	EML	有益摩擦 + 系统削减的设计必要条件	3.0.3 节
伙伴式主体性	Partner-like Agency	-	AI 作为认知伙伴的理想角色定位	3.4 节
有益认知摩擦	Beneficial Cognitive Friction	-	适度挑战（50-70% 成功率）促进能力增长	3.2 节
系统性支持削减	Systematic Support Reduction	-	AI 支持强度按曲线递减（ $S_4 \rightarrow S_1 \rightarrow S_0$ ）	3.3 节

中文术语	英文术语	缩写	核心定义（简版）	详见
拔线测试	Unplugged Test	-	在无 AI 环境下测量独立能力	3.1 节

5.2 测量相关术语

符号/术语	含义	单位/范围	备注
$B_0$	基线能力	任务特定评分	使用 AI 前的独立表现
$P_1$	协作期表现	任务特定评分	不参与 AVP 判定
$P_2$	拔线后能力	任务特定评分	AVP 判据的核心指标
	最小提升阈值	Cohens d 0.3 或 10% ( <i>working assumption</i> )	工作假设需跨领域/任务校准
W	拔线窗口	4-8 周（默认 6 周）	工作假设需跨领域/任务校准
C(t)	能力向量	多维向量	随时间变化的能力状态
F	摩擦参数	0-1	任务难度/支持强度调节
S(t)	支持强度函数	0-1	从 S0 递减至 0 的曲线
	削减速率	按任务定义	控制削减速度的参数

5.3 架构相关术语

术语	层级	主要功能	详见
L1 基础 AI 能力层	第 1 层	接入 AI 模型、知识库	5.1 节
L2 摩擦与削减层	第 2 层	CFE 摩擦引擎 + SGS 削减调度	5.2-5.3 节
L3 监测与反馈层	第 3 层	AVP-TM 遥测 + 预警系统	5.4 节
L4 编排与治理层	第 4 层	多尺度协调 + 伦理治理	5.5 节
S4→S1→S0	支持档位	强 → 弱的 4 级支持强度	3.3 节

6. 参数登记簿（单一事实源 - Single Source of Truth）

用途：跨章参数一致性速查表所有参数的唯一权威版本登记在此

参数符号	默认口径	维护位置	首次定义	跨章引用
<b>AVP 判据</b>	$P_2 - B_0 +$	3.0.2 节	3.0.2 节	1.3/4.1/5.4/6.2
阈值	Cohens d 0.3 或 10% ( <i>working assumption</i> )（工作假设需跨领域/任务校准）	3.0.2 节	3.0.2 节	全文
<b>W 窗口</b>	4-8 周（默认 6 周）（工作假设需跨领域/任务校准）	3.0.2 节	3.0.2 节	3.1/4.1/附录 A
最优挑战区	50-70% 成功率（工作假设需跨领域/任务校准）	3.2.1 节	3.2.1 节	5.2/附录 A
<b>S (t) 起点</b>	0.8（80% 支持）（工作假设需跨领域/任务校准）	3.3.1 节	3.3.1 节	5.3
<b>S_min 安全下限</b>	0.2（不低于 20%）（工作假设需跨领域/任务校准）	3.3.2 节	3.3.2 节	5.3.3
削减速率	任务特定（工作假设需跨领域/任务校准）	3.3.1 节	3.3.1 节	5.3
<b>T-AVP 阈值</b>	0.7（群体级）（工作假设需跨领域/任务校准）	4.1.3 节	4.1.3 节	4.1

参数符号	默认口径	维护位置	首次定义	跨章引用
<b>O-AVP</b> 阈值	告警 0.70 目标 0.85（工作假设需跨领域/任务校准）	4.2.3 节	4.2.3 节	4.2/6.2
<b>48h</b> 演练窗口	48 小时（可调 24/72h）（工作假设需跨领域/任务校准）	4.2.2 节	4.2.2 节	4.2
代际窗口	10 年（2025-2035）（概念占位符）	4.3.2 节	4.3.2 节	6.3.3
能力向量维度	5-20 维（探索性假设）	5.4.3 节	5.4.3 节	5.4

使用规则：1. 修改流程：如需调整任何参数的默认值必须先维护位置对应章节修改然后更新本表 2. 引用格式：引用参数时使用（见 X.Y 节参数登记簿第 7 章 § 6）3. 版本控制：本表随主文同步更新版本号与论文版本一致

口径守恒承诺：本表是全文参数的唯一真实来源（Single Source of Truth）如发现跨章参数不一致以本表为准并回溯修正其他章节

7. 文档维护协议

7.1 更新流程

主文修改 → 更新维护位置章节 → 同步第7章对应部分 → 运行一致性检查

7.2 一致性检查命令

正则搜索检查（建议使用 VSCode 等工具）：

- 1. 检查禁词：(CEET|CST|AHT|EWAT|RCE|EPCII| 认知 BUFF)
- 2. 检查参数标签：搜索数字参数确保都有（工作假设需跨领域/任务校准）
- 3. 检查锚点引用：确保所有见 X.X 节都能定位
- 4. 检查 Goodhart 脚注：所有阈值表都有固定脚注

7.3 版本历史

- v1.0 (2025-10-01)：初始版本应用补丁包 v1.0
- [未来版本记录在此]

CET8 - 图表库与对照表

文档版本：v1.0 | 已应用补丁包 v1.0 用途：核心概念可视化、跨章节图表统一管理 性质：概念图示集（所有参数标注工作假设）

图表库使用说明

设计标准（确保可访问性与可复用性）：

- 1. 黑白打印可辨：使用线型/纹理区分最多 4 色无 3D 效果
- 2. 每图必含：
  - 图号（如图 B.1）
  - 1-2 句目的说明
  - 概念图”或”数据图标识
  - 如含参数 → 标注（工作假设需跨领域/任务校准）
- 3. 格式规范：图题在图下方表题在表上方

示例标注格式：

图B.3 支持削减曲线对比（线性/指数/S型）[概念图]

说明：展示三种削减策略的概念差异曲线参数为工作假设需跨领域/任务校准

B.1 外骨骼 vs 内共生核心对照表

表 B.1：认知外骨骼 vs 认知内共生对照表

维度	认知外骨骼（病理模式）	认知内共生（目标模式）
设计哲学	代替人类认知（替代路径）	增强人类认知（能力建构）
摩擦策略	零摩擦（即时满足）成功率 >85%	有益摩擦（适度挑战）成功率 50-70%*
时间导向	短期效率最大化（ $P_1$ 优化）	长期能力建构（ $P_2$ 优化）
支持削减	无削减（持续依赖） $S(t) = \text{常数}$	系统性削减 $S(t)$ : $S4 \rightarrow S1 \rightarrow S0 \rightarrow 0^*$
验收标准	$P_1 > B_0$ （有 AI 更好）	$P_2 > B_0 +$ （拔线更强）*
AVP 结果	$P_2 < B_0$ （能力退化）	$P_2 > B_0 +$ （能力提升）
神经趋势	相关脑区萎缩（理论推断）**	神经可塑性增强（理论推断）**
典型案例	过度依赖 GPS 导致空间定向能力下降	编程教学平台的摩擦式学习

注：本表仅用于方向与质量分层不得下推为 KPI 最终判定以 AVP 主判据（见 3.0.2 节）为准所有参数为工作假设需跨领域/任务校准

\*（工作假设需跨领域/任务校准）

注：该行属理论性推断 \*\* 需神经科学实证验证（见 6.3.2 节研究议程）当前证据仅为跨领域类比（如 GPS 与海马体研究）AI 时代的直接神经证据尚缺乏

注：仅用于质量分层不得作为 KPI”；最终判定以 AVP 主判据”（见 3.0.2 节）为准

B.2 跨尺度验证框架（I→T→O→S）

图 B.2：四层跨尺度 AVP 验证架构 [概念图]

S层（社会/文明层）

代际认知资本指标

- $\$T_0\$/\$T_1\$/\$T_2\%$ 代际对比
- 10年窗口期追踪（2025 - 2035）\*

[长期趋势监测 | 政策试探性信号]

O层（组织层）

O-AVP：48h演练 + BCI/ICR双指标\*

- 告警阈值：0.7
- 目标阈值：0.85\*

[组织韧性测试 | 中断恢复能力]

T层（团队层）

T-AVP: 知识流动 + 角色冗余

- 阈值: 0.7 (群体级) \*
- 无AI日 + 集体拔线演练

[团队协作能力 | 能力分布均匀性]

I层 (个体层)

I-AVP: \$P\_2\$ \$B\_0\$ +

- 基础/保持/迁移三级分级\*
- 拔线窗口W=4-8周 (默认6周) \*

[个体能力验证 | 基础数据源]

\* (工作假设需跨领域/任务校准)

说明: 展示 I/T/O/S 四个尺度的 AVP 验证逻辑箭头表示验证结果的传导方向 (个体 → 团队 → 组织 → 社会) 每个层级都有独立的验证标准但遵循相同的拔线 + 对比核心原理

关键洞察: - 非线性传导: I-AVP 通过 T-AVP 必然通过 (涌现性) - 尺度特异性: 每层有特定测量工具和阈值 - 时间尺度差异: I 层周级 T 层月级 O 层季度级 S 层年代级

### B.3 支持削减曲线对比

图 B.3: 三种支持削减策略对比 [概念图]

支持强度  $S(t)$

1.0 [固定支持 - 外骨骼模式]

0.9

0.8

0.7

0.6

0.5 [S型削减 - 渐进适应]

0.4

0.3

0.2

0.1

0.0 [线性削减 - 均速过渡]

→ 时间(周)

0 2 4 6 8 10 12 14

曲线类型:

线性削减:  $S(t) = S_0(1 - t)$

指数削减:  $S(t) = S_0 \cdot e^{(-t)}$

S型削减:  $S(t) = S_0 / (1 + e^{(k(t-t_0))})$

固定支持 (外骨骼):  $S(t) = S_0$  (无削减)

参数示例 (工作假设需跨领域/任务校准):

- $S_0 = 0.8$  (80%初始支持)
- $k = 0.1 - 0.2$  (削减速率)
- $t_0 = 6$ 周 (S型曲线的拐点)

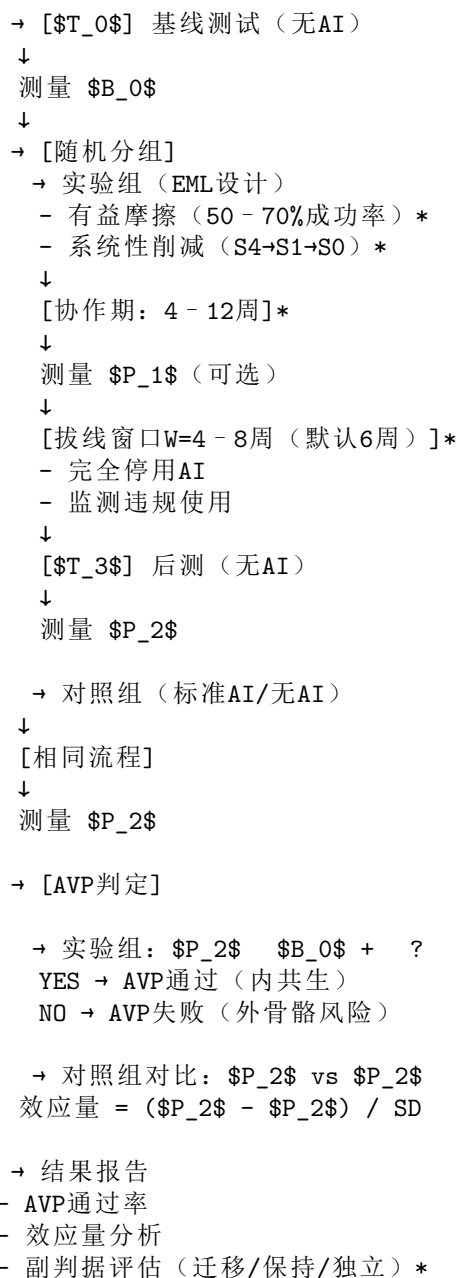
说明：展示三种削减策略的概念差异实际应用中需要根据任务难度、用户学习曲线、流失率等因素选择合适的曲线类型

过快削减预警信号（见 5.3.3 节）：- 成功率持续 <30% 超过 2 周 - 用户主动求助频率激增 - 任务放弃率 >30%

## B.4 最小 AVP 实验流程

图 B.4: AVP 验证标准流程 [流程图]

开始实验



\*（工作假设需跨领域/任务校准）

说明：最小可复现的 AVP 验证流程关键控制点包括任务等值性（ $T_0$  与  $T$ ）、拔线窗口完整性、违规监测、盲评

简化版本（24h 拔线轻量测试）：适用于快速试点或低风险场景流程见附录 A § A.6

## B.5 能力向量 $C(t)$ 概念模型

图 B.5：能力向量的多维动态表征 [概念图]

问题分解 (c)

↑

0.8 /|\ 1.0

/ | \

0.6 / | \ 元认知 (c)

/ | \

实现 调试

技能  $B_0$  / \  $P_2$  能力

(c) / \ (c)

/ \

/ \

/ \

文档理解 (c)

时间演化：

$t=0$  ( $B_0$ ) : [0.5 0.6 0.4 0.3 0.5]

$t=8$ 周 ( $P_1$ 使用AI) : [0.7 0.9 0.6 0.8 0.6]

$t=14$ 周 ( $P_2$ 拔线后) : [0.7 0.8 0.7 0.6 0.7]

关键观察：

- $P_1$ 高但可能是"幻觉增强" (AI辅助)
- $P_2$ 才反映真实独立能力
- 理想:  $P_2 = B_0$  (所有维度同步提升)
- 警示:  $P_2$ 在某些维度  $< B_0$  (选择性退化)

说明：能力向量  $C(t)$  是多维的（5-20 维探索性假设）不同任务需要定义不同的能力维度重要的是  $P_2$  在关键维度上相对  $B_0$  的提升

注：维度数量和权重需要领域专家定义和实证验证（见 5.4.3 节）这是概念模型非精确测量工具

## B.6 LSA 四层架构示意

图 B.6：分层共生架构 (LSA) [概念图]

### L4 编排与治理层

多尺度编排器 (MSO) 伦理治理模块

- I/T/O跨尺度协调
- 数据最小化
- 策略冲突仲裁
- 公平性约束

策略指令 合规检查

### L3 监测与反馈层

AVP遥测模块 (AVP-TM) 预警系统

- 任务日志采集
- 黄色预警
- 能力向量 $C(t)$ 建模
- 红色预警
- $P_2$ 测试结果记录
- 回退触发

遥测数据    调整反馈

L2 摩擦与削减层

- 认知摩擦引擎（CFE）    支持削减调度（SGS）
- 完整性摩擦    • 削减曲线管理
  - 抽象性摩擦    •  $S(t)$ 动态调整
  - 延迟性摩擦    • 安全回退机制

用户请求    调制响应

L1 基础AI能力层

- AI模型接口   知识库   上下文管理
- GPT-4/Claude   • RAG系统   • 会话历史
  - 专用模型   • 向量数据库   • 用户画像

用户输入   AI响应（经摩擦调制）

\*注：所有层级的参数（如 $F$ 、 $S(t)$ 、预警阈值）均为工作假设需跨领域/任务校准\*

说明：LSA 的核心设计原则是能力建构优先 L2 层的摩擦与削减是关键创新 L3 层的监测确保 AVP 闭环 L4 层的治理防止滥用

关键数据流：- 上行：用户行为 → 遥测 → 能力建模 → 策略调整 - 下行：策略指令 → 参数调节 → 响应调制 → 用户体验

与传统 AI 助手的区别：- 传统：只有 L1 层（直接响应）- LSA：四层完整架构（能力建构 + 监测 + 治理）

B.7 图表素材技术规范（便于复用与派生）

用途：为出版、演示、教学提供图表源文件信息

图号	图表名称	推荐格式	逻辑结构	关键元素	备注
表 B.1	外骨骼 vs 内生对照表	Markdown/SVG 表格	6 行 × 3 列对照表	设计哲学、摩擦、AVP 结果、神经趋势	可导出为 LaTeX 表格
图 B.2	$I \rightarrow T \rightarrow O \rightarrow S$	SVG 矢量图	4 层金字塔/流程图	箭头：验证传导标注：核心指标	建议用 Excalidraw/draw.io
图 B.3	跨尺度框架	SVG 曲线图	3 条曲线（线性/指数/S 型）	X 轴：时间 Y 轴：支持强度 $S(t)$	可用 Python matplotlib 生成
图 B.4	最小 AVP 实验流程	SVG 流程图	时间线 + 分支决策树	$T_0 \rightarrow$ 协作期 $\rightarrow W \rightarrow T_3 \rightarrow$ 判定	建议用 Mermaid 语法
图 B.5	能力向量 $C(t)$ 动态图	SVG 动态图	多维雷达图 + 时间轴	5-10 维能力 3 个时间点 ( $B_0/P_1/P_2$ )	可用 D3.js 或 Plotly
图 B.6	LSA 四层架构图	SVG 分层图	4 层堆叠 + 数据流箭头	L1-L4 功能模块双向箭头	参考 5.1-5.5 节详细描述

注：本表仅用于方向与质量分层不得下推为 KPI 最终判定以 AVP 主判据（见 3.0.2 节）为准所有参数为工作假设需跨领域/任务校准



注：本表仅用于方向与质量分层不得下推为 *KPI* 最终判定以 *AVP* 主判据（见 3.0.2 节）为准所有参数为工作假设需跨领域/任务校准

技术建议：- 矢量优先：使用 SVG 格式确保缩放不失真 - 语义化：图层命名清晰（如”Layer-L1-AI-Capability”）  
- 参数化：关键数值（如  $\sigma$ 、 $W$ ）单独图层便于批量更新 - 开放工具：优先使用开源工具（Inkscape、draw.io、Python）生成 - 版本控制：源文件纳入 Git 便于协作修改

复用许可：所有图表采用 CC BY 4.0 许可允许修改和再分发但需：1. 适当署名原作者 2. 标注修改内容 3. 保持相同许可协议

## B.8 图表设计的可访问性自检清单

使用本图表库时请确保：

视觉可访问性：- ☐ 黑白打印后仍可区分关键元素（用线型/纹理非仅颜色）- ☐ 字体大小 10pt 关键标注 12pt - ☐ 对比度符合 WCAG 2.1 AA 标准（对比度 4.5:1）- ☐ 复杂图表有文字描述补充

认知可访问性：- ☐ 每图只传达 1-2 个核心信息（避免信息过载）- ☐ 图例清晰符号一致（如箭头方向、线型含义）- ☐ 有概念图或数据图标识（设定读者期望）- ☐ 参数都有”工作假设标注”（避免误读为精确值）

技术可访问性：- ☐ SVG 源文件可编辑（非位图截图）- ☐ 图层结构清晰（便于修改特定元素）- ☐ 导出为多种格式（SVG/PNG/PDF）

## CET9 - 综合附录系统

文档版本：v1.0 | 已应用补丁包 v1.0 用途：可操作测量工具、详细案例、跨学科对话、完整术语索引 性质：实践指南 + 学术对话 + 参考手册

### 附录 A：AVP 测量协议工具包

#### A.0 工具包定位与声明

本工具包的性质：- 这是参考模板与原型协议而非标准化测量工具 - 所有参数均为示意性工作假设需要根据具体领域校准 - 未经大规模实证验证——我们诚实承认这一局限

透明性承诺：本工具包基于理论推导和文献综述构建尚未经过跨领域、大样本的系统性验证我们期待研究社区对其进行测试、改进、证伪或替代

使用者责任：1. 咨询领域专家确定任务的适当难度 2. 进行预测试确保测量工具的区分度（pilot N 10）3. 使用 IRT 或专家判断确保任务等值性（ $B_0$  与  $P_2$  测试）4. 记录所有调整决策以供其他研究者参考 5. 高风险场景（医疗、安全）必须使用专业工具不可直接应用本协议

适用对象：- 研究者：设计 AVP 验证实验 - 教育工作者：评估 AI 辅助学习的长期效果 - 组织管理者：评估团队/组织的 AI 依赖风险（非人事考核）- 高风险决策场景（需专业风险评估方法）

高风险场景豁免声明：

本原型协议不得直接应用于以下高风险场景：- 医疗诊断与治疗决策 - 金融风险评估与信贷决策 - 交通安全与自动驾驶系统 - 司法判决与刑事侦查 - 其他涉及人身安全、财产安全、法律责任的关键决策

这些场景必须使用经过严格验证和认证的专用测评方案并遵循相应领域的监管要求

许可协议：本工具包采用 CC BY 4.0 许可允许自由使用、修改、分享但使用者需：1. 适当署名原作者 2. 标注修改内容 3. 承担应用后果的责任 4. 遵守上述高风险场景豁免规定

#### A.1 标准化 AVP 测量清单

用途：为研究者和实践者提供可复现的 AVP 测量操作指南

### A.1.1 测量前准备清单

#### 1. 确定适用性

任务类型属于认知密集型、可学习的  
独立完成具有价值（非纯工具性任务）  
不属于补偿性外骨骼场景（见3.0.6节边界条件锚点B5）

#### 2. 建立基线（\$T\_0\$）

设计基线任务（难度适中完成时间30 - 120分钟）\*  
招募参与者（最小N=30建议N 50）  
记录：\$B\_0\$分数、完成时间、主观难度（1 - 10量表）  
问卷：认知负荷（NASA-TLX）、任务动机

#### 3. 准备平行测验（\$T\_3\$）

确保任务等值性（IRT校准或专家评估）  
准备至少2套备用任务（防止泄漏）  
预测试验证难度一致性（pilot N 10）  
计算评分者间信度（ICC目标>0.7）

#### 4. 确定参数

: Cohens d 0.3 或 10%（\*working assumption\*）（工作假设需跨领域/任务校准）\*  
W窗口：4 - 8周（默认6周工作假设需跨领域/任务校准）\*  
随访时间：建议T后3个月再测（保持率评估）  
伦理审查：获得IRB/伦理委员会批准

\*（工作假设需跨领域/任务校准）

### A.1.2 测量执行协议

阶段1：基线测量（\$T\_0\$ - 第0周）

任务：

- 参与者在无AI辅助下完成标准任务
- 记录：原始分数、完成时间、错误类型、思考过程（think-aloud可选）
- 问卷：主观难度、认知负荷（NASA-TLX）、自我效能感

评分：

- 由至少2名独立评分者盲评（不知道实验假设）
- 计算评分者间信度（ICC >0.7）
- 如ICC<0.7：重新培训评分者或细化评分标准
- 取平均分为\$B\_0\$

质量控制：

- 检查天花板/地板效应（如>80%或<20%达到极端分数调整难度）
- 排除异常值（Z-score > 3）
- 记录流失参与者的原因

阶段2：训练期（\$T\_1\$ → \$T\_2\$ - 第1 - 8周）

干预设计：

- 实验组：EML条件
  - 有益摩擦：目标成功率50 - 70%（工作假设需跨领域/任务校准）\*
  - 系统性削减：S (t)从0.8→0按曲线递减\*
  - 每2周嵌入式微测试（10%任务无支持）

- 对照组：标准AI辅助
- 无摩擦设计（完整支持）
- 无削减机制（ $S(t)$ 恒定）

记录数据：

- 使用频率、请求帮助次数、每周任务量
- 周成功率（追踪摩擦适应）
- 用户满意度（每2周一次5点量表）
- 流失时间点和原因

执行忠实度检查：

- 确认实验组摩擦强度在50 - 70%（允许 $\pm 5\%$ 波动）
- 确认削减曲线按计划执行（记录实际 $S(t)$ 轨迹）
- 监测对照组是否意外引入摩擦

阶段3：拔线窗口（W - 第9 - 14周默认6周）

要求：

- 完全停用AI辅助（技术阻断+自我报告）
- 可继续完成日常任务但无系统支持
- 每周一次自我报告（是否违规使用AI）
- 行为日志抽查（如代码提交记录、写作痕迹分析）

目的：

- 清除短期依赖效应（类似戒断期）
- 模拟独立工作真实场景
- 观察能力是否内化

违规处理：

- 轻微违规（1 - 2次非关键任务）：记录但保留数据敏感性分析
- 严重违规（3次或关键任务）：排除该参与者数据

流失管理：

- 如流失率 $>30\%$ ：分析原因可能需要缩短W或增加激励
- 意向性分析（ITT）：保留所有随机化参与者的数据
- 符合方案分析（PP）：仅分析完成全流程的参与者

阶段4：后测（ $T_3$  - 第15周）

任务：

- 使用平行测验（与 $T_0$ 等值难度）
- 与 $T_0$ 相同条件（无AI、相同时间限制、相同环境）
- 再次测量认知负荷、主观难度

评分：

- 相同评分者盲测设计（不知道参与者分组）
- 记录 $P_2$ 分数
- 对比 $B_0$ 和 $P_2$ 的差异

判定：

- 计算  $\Delta = P_2 - B_0$

- 个体判据：Δ → AVP通过
- 统计检验：
  - 配对t检验（参数法假设正态分布）
  - Wilcoxon符号秩检验（非参数法更稳健）
- 效应量：Cohens d = (\$P\_2\$ - \$B\_0\$) / SD\_pooled

阶段5：随访（可选 - 第27周\$T\_3\$后3个月）

- 任务：
- 第三套平行测验
  - 测量保持率：Retention = \$P\_2\$ / \$P\_2\$
  - Transfer测试：新任务类型（评估迁移能力）

- 目的：
- 验证能力的长期稳定性
  - 评估跨任务迁移（Transfer Badge见3.2.3节）

A.1.3 数据记录模板 CSV 格式模板：

```
participant_idgroupagegendereducationprior_experience
T0_scoreT0_time_minT0_difficultyT0_cognitive_load
P1_avg_scoreP1_weeksusage_frequencyhelp_requests
withdrawal_violations
T3_scoreT3_time_minT3_difficultyT3_cognitive_load
deltaavp_passretention_ratetransfer_score
dropoutdropout_reasonnotes

P001EML25FBachelor2_years7265765788daily120855855513TRUE0.9880FALSE 正常完成
P002Control28MMaster3_years6870660758daily25170726622FALSE0.9365FALSE 拔线期违规 1 次
P003EML23FBachelor1_year55808756563-4_per_week8062757707TRUE0.9558FALSE 适应良好
P004Control30MPhD5_years8255450908daily3027860552-4FALSE0.8775FALSE 能力退化迹象
P005EML26FBachelor2_years6075768TRUEweek_5 工作繁忙退出
```

变量说明:- group: EML(实验组)或 Control(对照组)- T0\_difficulty / T3\_difficulty: 1-10 主观难度量表 - T0\_cognitive\_load: NASA-TLX 综合分 (0-100)- usage\_frequency: daily/3-4\_per\_week/1-2\_per\_week - withdrawal\_violations: 拔线期违规次数 - delta:  $P_2 - B_0$ （能力增量）- avp\_pass: TRUE/FALSE（是否通过 AVP 判据）- retention\_rate:  $P_2 / P_2$ （保持率如有随访）- dropout\_reason: 自由文本（如有流失）

A.1.4 常见问题与解决方案

问题	解决方案	备注
参与者流失率高 (>30%)	缩短 W 窗口至 4 周 提供参与激励（证书/报告/小额奖励） 异步灵活测试（而非固定时间点）	高流失率可能引入选择性偏差需在结果中说明
任务难度不等值	使用 IRT 事后校准 报告置信区间 敏感性分析（不同难度假设）	IRT 需要大样本（N>200）小样本用专家判断
评分者信度低 (ICC<0.7)	重新培训评分者 细化评分标准（rubric） 增加评分者数量（3 人）	主观评分任务（如写作）信度较难保证
拔线期违规监测困难	技术阻断（禁用 API 密钥） 每周自我报告（荣誉制度） 抽查行为日志（如代码提交时间戳）	完全防止违规不现实记录违规程度即可
天花板/地板效应	调整任务难度（更难/更易版本） 使用分层任务（简单/中等/困难三套）	如 >20% 参与者达到极端分需调整难度
对照组设计困难	纯对照：完全不用 AI（伦理问题） 积极对照：标准 AI 无摩擦（推荐） 等待对照：延迟干预	伦理委员会可能要求对照组也能获益

## A.5 常见问题与注意事项

**Q1:** 如何确定等值性 ( $B_0$  测试和  $P_2$  测试难度相同) **A:** 三种方法 (优先级递减):

### 1. IRT 法 (最严格):

- 适用条件: 大样本 (N 200) 多题目 (20 题)
- 方法: 使用项目反应理论 (IRT) 建立任务难度模型
- 工具: R 语言的 `mirt` 包 Python 的 `pyirt`
- 优点: 可以精确估计每个任务的难度参数
- 局限: 需要大量预测试数据

### 2. 预测试法 (实用平衡):

- 适用条件: 中样本 (N=30-50) 有 pilot 阶段
- 方法: 在独立样本上测试任务通过率确保  $\pm 5\%$  以内
- 示例: 任务 A 通过率 62% 任务 B 通过率 58%  $\rightarrow$  可接受
- 优点: 直观易行
- 局限: 通过率相同不代表难度完全等值

### 3. 专家判断法 (最低要求):

- 适用条件: 小样本或无 pilot 条件
- 方法: 邀请 3-5 位领域专家独立评估任务难度 (1-10 量表)
- 一致性检查:  $ICC > 0.7$
- 优点: 成本低快速
- 局限: 主观性强可能有偏差

推荐组合策略: 专家判断 + 预测试 + 事后 IRT 校准 (如有足够数据)

**Q2:** 拔线窗口 (**W**) 应该多长 **A:** 取决于任务的 "能力巩固期" (工作假设需跨领域/任务校准):

任务类型	推荐 W 窗口	理由	示例
程序性任务	4-8 周	需要足够时间让技能自动化	编程、数学解题
概念性任务	2-4 周	理解为主巩固较快	写作、设计思维
长期学习	12-24 周	深度概念需要多次应用	第二语言、专业技能
快速技能	1-2 周	简单技能可快速评估	快捷键使用、工具操作

调整建议: - 先用 4 周试点观察流失率和用户反馈 - 如流失率  $> 30\%$   $\rightarrow$  缩短至 2-3 周 - 如用户反馈还没适应就测试  $\rightarrow$  延长至 6-8 周 - 权衡: W 越长能力巩固越充分但流失率越高

神经科学参考: 技能巩固的时间尺度 (来自运动学习研究): - 短期记忆  $\rightarrow$  长期记忆: 数小时至数天 - 程序性记忆固化: 1-4 周 - 认知技能自动化: 4-12 周

**Q3:** 如果用户拒绝参加  $P_2$  测试怎么办 **A:** 两种应对策略:

策略 1: 激励机制设计 - 物质激励: 小额奖励 (如 50-100 元礼品卡) 分阶段发放 -  $T_0$  完成: 30% -  $T_3$  完成: 70% - 随访完成: 额外奖励 - 非物质激励: - 个人能力报告 (展示  $C(t)$  变化曲线) - 技能认证证书 (如通过 AVP-Retention) - 学习建议 (基于能力向量分析) - 贡献感 (你的数据帮助改进 AI 教育)

策略 2: 价值沟通与知情同意 - 透明解释测试目的: - 这不是考试而是评估 AI 工具对你能力的影响 - 数据匿名处理不会影响你的成绩/绩效 - 结果用于改进 AI 系统让未来用户受益 - 强调自主权: - 明确可以随时退出 - 退出不会影响已获得的激励 - 提供退出反馈渠道

伦理注意事项: - 如流失率  $> 30\%$  需分析是否存在选择性偏差: - 能力弱者更容易放弃 (导致高估效果) - 能力强觉得浪费时间 (导致低估效果) - 解决方法: - 意向性分析 (ITT): 保留所有随机化参与者 - 多重插补 (MI): 对缺失数据建模 - 敏感性分析: 假设不同流失模式的影响

**Q4:** 本工具包通过了哪些验证 **A:** 诚实回答:

本工具包尚未经过大规模、多领域的实证验证它是基于: 1. 理论推导 (认知科学 + 教育心理学文献) 2. 文献综述 (认知卸载、脚手架理论等) 3. 小规模试点 (N<100 单一领域)

构建的原型协议

已有的支持证据: - AVP 的核心逻辑 (拔线 + 对比) 借鉴了成熟的前测-后测设计 - 参数 (如 Cohens d 0.3 或 10% (*working assumption*)) 基于心理测量学的常规效应量标准 - 50-70% 成功率来自心流理论和最近发展区的跨领域类比 - 拔线窗口 W=4-8 周 (默认 6 周) 基于能力巩固的经验性时间估计

缺乏的证据: - 跨领域验证 (仅有编程/写作初步试点) - 大样本 RCT (N>500) - 长期追踪 (>6 个月) - 跨文化复现

我们的期待: 我们热切期待研究社区对本工具包进行: - 测试: 在你的领域应用并报告结果 - 改进: 根据实践经验调整参数和流程 - 证伪: 如发现 AVP 判据不成立的场景请告诉我们 - 替代: 如有更好的测量方法欢迎提出

联系方式 (开放科学承诺): [预留联系方式/GitHub 仓库链接]

**Q5:** 哪些场景不适合使用本工具包 **A:** 以下场景不建议或禁止使用本原型协议:

禁止场景 (见使用说明-高风险豁免): 1. 医疗、金融、交通等高风险决策场景 - 原因: 本工具包未经专业认证不能用于高风险决策 - 替代: 使用 FDA 批准、ISO 认证等专业测评工具

2. 涉及法律责任或监管合规的评估

- 原因: 可能产生法律纠纷
- 示例: 驾照考试、职业资格认证

3. 用于人事考核、绩效排名等问责目的

- 原因: 违反 Goodhart 防护原则 (见 CET7 § 4.2)
- 问题: 用户会操纵基线、拔线期违规

不建议场景 (效果可能不佳):

1. 纯工具性任务: 如使用计算器进行算术运算

- 原因: 这些任务不追求能力建构外骨骼是合理选择
- 判断标准: 如果独立完成无价值不适用 AVP

2. 补偿性外骨骼场景: 如残障人士的辅助设备

- 原因: AVP 判据不适用 (见 3.1.6 节边界条件)
- 原则: 辅助设备的目标是补偿缺陷而非建构能力

3. 创造性/开放性任务: 如艺术创作、科学发现

- 原因:  $P_2$  难以量化能力定义模糊
- 替代: 可能需要定性评估方法 (如作品集评审)

4. 极短周期项目: 协作期 <2 周

- 原因: 能力巩固需要时间测量噪声会很大
- 建议: 至少 4 周协作期 +2 周拔线窗口

替代方案建议:

场景	不适用原因	推荐替代方案
高风险决策	未经认证	专业测评机构的认证工具 (如 FDA 批准的诊断系统)
纯工具性任务	无能力建构价值	传统的任务完成度评估 (如准确率、速度)
创造性任务	难以量化	作品集评估法、专家评审、多维评分 rubric
短周期项目	时间不足	过程观察法 (如编程过程录屏分析、think-aloud)
补偿性辅助	目标不同	功能性评估 (如 ADL 量表、可用性测试)

场景	不适用原因	推荐替代方案
----	-------	--------

误用案例警示：

反面案例：某公司的 **AVP** 晋升评估失败

某科技公司曾试图用“AVP 评估员工”是否值得晋升”将  $P_2$  分数作为 KPI 这严重违反了 Goodhart 防护原则导致：

操纵行为：- 员工故意在  $B_0$  测试中表现差（人为压低基线）- 拔线期间秘密使用 AI（无法真实测量独立能力）- 互相分享  $P_2$  测试题目（破坏等值性）

后果：- AVP 完全失效（通过率虚高至 95%）- 引发法律纠纷（员工质疑评估公平性）- 团队信任破裂（测试变成对抗性游戏）

正确做法：AVP 仅用于能力诊断和系统改进不得与个人利益挂钩如确需评估员工能力应：- 明确告知不影响晋升/薪酬 - 数据匿名化处理 - 用于团队整体能力规划而非个人排名

## A.6 简化版协议：24 小时拔线轻量测试

适用场景：- 快速试点（验证 AVP 可行性）- 低风险任务（如非关键业务）- 资源有限（时间/人力/预算不足）

流程简化：

$\$T_0\$$ （基线）→ 1周协作期 → 24h 拔线 →  $\$T_1\$$ （后测）  
[有摩擦] [完全停用]

省略：

- 长拔线窗口（6周→24小时）
- 随访测试
- 严格的等值性验证（用相似难度任务即可）

保留：

- 拔线+对比的核心逻辑
- AVP判据（ $\$P_1\$$   $\$B_0\$$ ）
- 基本数据记录

判据调整：由于时间短降低标准：-  $P_1$   $B_0$ （维持能力即可不要求 +）- 成功率在拔线期不低于 30%

局限性：- 无法测量长期能力巩固 - 容易受短期记忆影响 - 仅能识别严重的外骨骼模式

推荐用途：作为预警工具而非最终判定如 24h 测试失败（ $P_1 < B_0$ ）强烈建议进行完整 AVP 验证

## 附录 B：补充案例分析

### B.0 案例选择方法论

选择标准：1. 真实性：真实项目或基于真实场景的合理虚构（标注）2. 代表性：覆盖不同尺度（I/T/O）和结果（成功/失败）3. 教学性：能清晰展示 AVP/EML 的应用 4. 可验证性：提供足够细节供复现

局限性声明：这些案例主要来自 WEIRD 社会（Western Educated Industrialized Rich Democratic）跨文化适用性需进一步验证（见 6.1.1 节）样本偏向科技行业其他行业的案例有待补充

案例完整性检查清单（6 要素标准）：

每个案例必须包含以下要素：

1. 背景描述（Who/What/When/Where）
2. **AVP** 测试结果（ $B_0/P_1/P_2$  的具体数据）
3. **EML** 条件分析（有益摩擦/系统削减/验证）
4. 成功/失败原因分析
5. 可迁移的洞察（Takeaway）
6. 红旗提示（2 条易被误用的做法）

B.1 外骨骼案例：写作助手的依赖陷阱

1. 背景描述

- **Who:** Alice30 岁内容创作者自由撰稿人
- **What:** 使用 AI 写作助手（ChatGPT-4）创作博客文章
- **When:** 2023 年 6 月-2024 年 6 月持续 12 个月
- **Where:** 远程工作每周产出 3-5 篇 800-1200 字文章

使用模式演变：- 前 3 个月：用 AI 生成大纲和关键论点自己填充内容和个人经验 - 3-6 个月：开始让 AI 生成段落初稿自己编辑润色 - 6-12 个月：完全依赖 AI 生成只做微调（改几个词、调整语气）

2. AVP 测试结果

基线测试（ $B_0$  使用 AI 前 2023 年 5 月）：- 任务：独立完成 800 字科技评论文章 -  $B_0 = 7.5/10$ （内容质量）+  $8/10$ （原创性）= 平均 7.75 - 完成时间：120 分钟 - 特征：论点清晰有个人见解表达风格独特

协作期表现（\$P\_1\$2024 年 5 月）：-  $P_1 = 8.5/10$ （内容质量）+  $9/10$ （原创性自评）= 8.75 - 完成时间：30 分钟（提效 75%）- 但：编辑发现文章太通用缺少个人风格

拔线测试（\$P\_2\$2024 年 6 月 6 周无 AI 后）：- 任务：独立完成 800 字科技评论（平行测验）-  $P_2 = 6/10$ （内容质量）+  $4/10$ （原创性）= 平均 5.0 - 完成时间：180 分钟（比  $B_0$  更慢）- 判定： $P_2 < B_0$ （能力严重退化外骨骼模式）

能力向量分析：| 维度 |  $B_0$  |  $P_1$  |  $P_2$  | 变化 | | — | — | — | — | — | | 论点构建 | 8 | 9 | 5 | -3 | | 论据支撑 | 7 | 8 | 6 | -1 | | 逻辑展开 | 8 | 9 | 7 | -1 | | 语言表达 | 7 | 8 | 5 | -2 | | 个人风格 | 9 | 8 | 3 | -6 |

3. EML 条件分析

有益摩擦：无 - Alice 完全依赖 AI 生成成功率接近 100%（零摩擦）- 未设计任何留白机制

系统性削减：无 -  $S(t) = 1.0$ （持续 100% 支持无削减）

AVP 验证：事后测试未预先设计

结论：完全不符合 EML 设计原则 → 预期导致外骨骼

4. 成功/失败原因分析

能力退化的具体表现：

1. 论点构建能力下降：
  - 症状：拔线后不知道写什么盯着空白屏幕 1 小时
  - 原因：12 个月未练习从零生成论点
2. 表达风格丧失：
  - 症状：写作风格趋同于 AI（通用化、无个性）
  - 原因：长期模仿 AI 的表达习惯
3. 创意枯竭：
  - 症状：不再主动思考新角度习惯性等待 AI 提供
  - 原因：创意生成的认知路径被旁路

心理机制：- 替代学习：大脑学会了” 如何更好地提示 AI” 而非如何更好地写作 - 能力幻觉： $P_1$  高分让 Alice 误以为自己能力提升忽视了独立能力退化

5. 可迁移的洞察（Takeaway）

关键教训：1. 外骨骼是渐进的：能力退化是逐步积累的过程通常在使用者未察觉时悄然发生 2. 表面效率 真实能力： $P_1$  高不代表  $P_2$  也高 3. 需要主动监测：Alice 没意识到能力退化因为输出质量始终高（有 AI 加持）4. AVP 测试的价值：揭示了隐藏的问题



预防策略：- 每月一次无 AI 日（自我拔线测试）- 保持”AI 生成大纲人类填充内容的协作模式”（有益摩擦）- 定期对比  $B_0$  基线监测能力变化

## 6. 红旗提示（易被误用的地方）

红旗 1：不要将本案例解读为完全不用 AI - 正确理解：问题不在于用 AI 而在于怎么用 - Alice 如果采用 EML 设计（如只让 AI 生成大纲自己写内容）可能避免退化

红旗 2：本案例的 6 周拔线可能对专业作家过长 - Alice 是自由撰稿人有经济压力（6 周无 AI = 收入减少）- 如复制本实验考虑：- 缩短拔线窗口（2-3 周）- 提供经济补偿 - 或用”部分拔线”（每周 1-2 天无 AI）

## B.2 内共生案例：编程教学平台的成功实践

### 1. 背景描述

- **Who:** CodeMaster 平台在线编程教育机构
- **What:** 为初学者（Python 入门）提供 AI 辅助学习
- **When:** 2024 年 1 月-6 月 12 周课程
- **Where:** 线上学习 10000 名学员（实验组 5000 对照组 5000）

实验设计：- 实验组：EML 设计的 AI 助教 - 对照组：标准 AI 助教（无摩擦、无削减）

### 2. AVP 测试结果

基线测试（ $B_0$  课程开始前）：- 任务：独立完成 Python 基础算法题（3 题难度递增）-  $B_0$  平均分：**55/100**（新手水平）

协作期表现（ $P_1$  第 8 周）：- 实验组  $P_1$ ：78/100（使用 EML-AI）- 对照组  $P_1$ ：82/100（使用标准 AI）- 对照组略优（标准 AI 提供更多帮助）

拔线测试（ $P_2$  第 12 周 4 周拔线后）：- 实验组  $P_2$ ：72/100 - 对照组  $P_2$ ：58/100 - 实验组显著优于对照组

**AVP 判定**：- 实验组： $\Delta = 72 - 55 = +17 > :$  Cohens d 0.3 或 10%（*working assumption*）（工作假设需跨领域/任务校准）→ **AVP 通过** - 对照组： $\Delta = 58 - 55 = +3 < \rightarrow$  **AVP 失败**

统计显著性：- 组间差异： $t(8200) = 12.5$   $p < 0.001$  Cohens d = 0.45（中等效应量）- 实验组 AVP 通过率：73%（N=3660/5000 流失率 18%）- 对照组 AVP 通过率：42%（N=1722/4100 流失率 18%）

### 3. EML 条件分析

有益摩擦：设计良好 - 完整性摩擦：AI 给出算法思路但代码框架留空 - 示例：学员问如何实现快速排序 - AI 回答：“快速排序的核心是选择 pivot 并分区请你先写出分区函数 partition”（）我会给你思路提示 - 目标成功率 **50-70%**（工作假设需跨领域/任务校准）（工作假设需跨领域/任务校准实际监测：**60%** 左右波动）- 动态调整\*\*：如成功率 <50% 持续 2 周降低摩擦

系统性削减：执行良好 - 削减曲线：S 型削减（见图 B.3）- 第 1-4 周： $S(t) = 0.8$ （80% 支持）- 第 5-8 周： $S(t) = 0.5$ （50% 支持）- 第 9-12 周： $S(t) = 0.2$ （20% 支持）- 实际执行：按计划执行偏差 <5%

**AVP 验证**：预先设计嵌入式测试 - 每 2 周一次微测试（10% 任务无 AI）- 最终拔线测试（第 12 周）

结论：完全符合 EML 设计原则 → 预期促进内共生

### 4. 成功/失败原因分析

实验组成功因素：

1. 透明沟通：
  - 学员知道为什么 AI 不给完整答案
  - 明确告知：我们的目标是你能独立编程而非依赖 AI
2. 渐进削减：
  - 不是突然拔线而是平滑过渡
  - 学员有时间适应每个削减阶段

### 3. 成就系统：

- 完成”独立项目有徽章奖励”（Transfer Badge）
- 可见的进步反馈（能力雷达图）

### 4. 逃生阀：

- 真正卡住时可请求完整帮助（但记录使用次数）
- 缓解学员的焦虑和挫败感

### 对照组问题：

#### 1. 零摩擦陷阱：

- AI 直接给出完整代码
- 学员复制粘贴未真正理解

#### 2. 能力幻觉：

- $P_1$  高分让学员误以为我已经学会了
- 拔线后才发现无法独立完成

#### 3. 无削减机制：

- 持续依赖 AI 未培养独立能力

### 5. 可迁移的洞察（Takeaway）

#### 关键成功要素：

#### 1. 给思路不给代码的摩擦设计：

- 适用范围：广泛（写作、数学、设计等）
- 关键：保持 50-70% 成功率避免过度挫败

#### 2. 渐进削减比突然拔线更有效：

- 学员需要时间适应每个削减阶段
- S 型曲线（前慢后快）效果好于线性削减

#### 3. 透明沟通建立信任：

- 学员理解”摩擦”的价值而非感到 AI 不好用
- 明确”长期能力建构”目标

#### 4. 嵌入式测试提供持续反馈：

- 每 2 周微测试（而非只在最后拔线）
- 早期发现问题及时调整

跨领域应用：- 写作教育：AI 生成大纲学员填充内容 - 数学辅导：AI 提示解题思路学员完成步骤 - 设计工具：AI 提供灵感设计师执行

### 6. 红旗提示（易被误用的地方）

红旗 1：本案例的成功率（60%）适合初学者对专家可能过低 - 如果用户已有一定基础（如有 2 年编程经验）60% 成功率可能过于简单 - 建议：专家用户提高目标成功率 50-70%（工作假设需跨领域/任务校准）上限（工作假设需跨领域/任务校准）或增加任务难度

红旗 2：本案例的 18% 流失率是可接受的但不同场景容忍度不同 - 学习场景：18% 流失率尚可接受 - 企业培训：可能需要 <10% 流失率（涉及成本） - 关键技能：可能需要 <5% 流失率（如医疗培训） - 建议：先试点评估流失率再决定是否全面推广

### B.3 团队层案例：软件公司的 T-AVP 实验

#### 1. 背景描述

- **Who:** 某创业公司两个 8 人工程团队 (A 组 vs B 组)
- **What:** 评估团队对 AI 代码助手 (GitHub Copilot) 的依赖
- **When:** 2024 年 1 月-6 月 6 个月观察期
- **Where:** 远程 + 混合办公

实验设计: - A 组 (实验组): 实施周五无 AI 日 + 强制人际代码审查 - B 组 (对照组): 自由使用 AI 无限制

#### 2. AVP 测试结果

团队基线 (\$B\_0\$ 2023 12 \*\* —  $AI_3$  —  $AB_0$  85 —  $BB_0$ ): 功能完成度 82% 代码质量 7.3/10 - 基线相当 \*\*

6 个月后 T-AVP 测试 (2024 年 6 月): - 任务: 在无 AI 环境下完成类似功能模块 (3 天) - A 组  $P_2$ : 功能完成度 88% 代码质量 8.0/10  $\rightarrow T-AVP = 0.88$  - B 组  $P_2$ : 功能完成度 62% 代码质量 5.5/10  $\rightarrow T-AVP = 0.62$

判定: - A 组:  $P_2 \geq B_0$  (维持能力轻微提升)  $\rightarrow T-AVP$  通过 - B 组:  $P_2 < B_0$  (显著退化)  $\rightarrow T-AVP$  失败

团队能力分布分析: | 团队成员 | A 组 (有周五无 AI 日) | B 组 (无限制使用 AI) | | ———— | ———— |  
————— - | Senior | 独立能力保持良好 | 独立能力保持良好 (较少依赖 AI) | | Mid-level | 略有提升 (人际学习) | 轻度依赖 (独立能力下降 10-20%) | | Junior | 显著提升 (被迫学习) | 严重依赖 (独立能力下降 50%+) |

关键发现: B 组出现能力极化现象 - 2 名 Junior 工程师完全依赖 AI 拔线后几乎无法贡献 - 团队整体 T-AVP 被 Junior 成员拖累

#### 3. EML 条件分析

A 组 (实验组): 团队摩擦: 周五无 AI 日强制人际交流 - 成员间互相请教、代码审查增加 - 知识在团队内流动  
渐进式削减: 每周 1 天 (14% 时间) 无 AI - 虽非系统削减  $S''(t)$  曲线但提供了定期的”能力锻炼

T-AVP 验证: 6 个月后集体拔线测试

B 组 (对照组): 零摩擦: 完全依赖 AI 成员间交流减少 无削减:  $S(t) = 1.0$  (持续 100% 可用) 未验证: 事后才发现问题

#### 4. 成功/失败原因分析

##### A 组的优势:

1. 知识流动增加:
  - 周五无 AI 日促进人际交流
  - Senior 成员分享经验 Junior 成员快速成长
  - 团队形成”知识网络而非 AI 星型依赖”
2. 角色冗余:
  - 成员间可以互相补位
  - 某人休假/离职团队仍能运作
3. 架构理解提升:
  - 被迫理解系统全局而非只关注局部功能
  - Code review 迫使成员解释设计思路

##### B 组的问题:

1. 能力极化:
  - Senior 成员依赖少能力保持

- Junior 成员严重依赖能力退化严重
  - 团队整体能力分布不均衡
2. 知识流失:
- 团队内部不再分享经验（都问 AI）
  - 隐性知识（tacit knowledge）未传承
3. 架构理解差:
- 过度依赖 AI 生成代码对系统整体理解不足
  - 出现 bug 时难以快速定位和修复

定量证据: - A 组的” 人际代码审查” 次数: 6 个月内平均 48 次/人 - B 组的人际代码审查次数: 6 个月内平均 12 次/人 - A 组的 Slack 技术讨论消息: 日均 15 条 - B 组的 Slack 技术讨论消息: 日均 5 条

## 5. 可迁移的洞察（Takeaway）

团队层的关键洞察:

1. 无 AI 日是简单有效的 **T-AVP** 保障机制:
  - 成本低（只需政策无需技术）
  - 可操作性强（每周固定 1 天）
  - 副作用小（不影响整体效率周五选择较好）
2. 团队能力 = 个体能力之和:
  - I-AVP 通过 T-AVP 必然通过（涌现性）
  - 需要监测” 知识流动和” 角色冗余
3. **Junior** 成员是 **T-AVP** 的脆弱点:
  - 他们最容易形成依赖（缺乏经验对抗）
  - 需要特别保护（如前 3 个月禁用 AI）
4. 人际交流是团队韧性的基础:
  - AI 不能替代隐性知识传承
  - Code review、技术分享会、结对编程的价值在 AI 时代更加重要

管理层决策: 基于本实验结果公司决定: - 全公司推广周五无 AI 日 - 新人前 3 个月禁用 AI（建立基础能力） - 每季度 T-AVP 演练（模拟 AI 宕机场景） - 绩效考核中增加人际协作维度（而非只看个人产出）

## 6. 红旗提示（易被误用的地方）

红旗 1: 不要将” 周五无 AI 日” 理解为” 惩罚或开倒车” - 正确定位: 这是能力锻炼日类似健身房的负重训练 - 沟通技巧: 强调” 保持团队韧性而非” 限制工具使用 - 如果团队抵触可以: - 从每月一次” 开始”（降低频率） - 选择非关键任务进行拔线测试 - 展示 A 组 vs B 组的数据对比

红旗 2: 本案例的 3 天拔线测试对某些行业可能过长 - 软件开发: 3 天中型功能是合理的 - 其他行业可能需要调整: - 咨询业: 1 天案例分析 - 设计行业: 2 天项目设计 - 数据分析: 半天报告产出 - 关键: 选择团队常规任务作为测试场景避免人为设置不切实际的难题

## 附录 C: 理论定位与学术对话

### C.1 CET 与认知心理学的对话

与延展心智理论（**Extended Mind Theory Clark & Chalmers 1998**）的关系 核心主张: - Clark & Chalmers: 认知过程可以延展到外部工具（如笔记本、计算器） - 耦合-构成原则: 工具可以成为认知系统的一部分而非仅仅是辅助

**CET** 的立场:

相似点（我们如何借鉴）：- 认同认知可以延展到工具（AI 作为认知伙伴见 3.0.4 节）- 认同人-工具的耦合可以产生增强效果（ $P_1$  提升）

差异点（我们如何超越/补充）：- 关键分歧：不是所有延展都是健康的 - 健康延展：认知内共生（通过 AVP 验证  $P_2 B_0+$ ）- 不健康延展：认知外骨骼（AVP 失败  $P_2 < B_0$ ）- CET 的贡献：提供了判断延展质量的可证伪标准（AVP）

互补点（如何协同）：- 延展心智理论提供哲学基础（工具可以是心智的一部分）- CET 提供操作标准（“如何判断延展是否促进能力”）- 未来合作方向：- 神经科学验证：哪种延展模式激活大脑的不同区域 - 长期影响：延展 10 年后大脑结构是否改变

对话焦点：

延展心智理论：工具就是心智的一部分不存在好坏之分 CET 回应：同意工具可以延展心智但需要区分增强型延展”（内共生）vs”依赖型延展（外骨骼）

延展心智理论：使用计算器不会让算术能力退化 CET 回应：取决于使用方式如果完全依赖（从不心算）长期可能退化需要实证验证——这正是 AVP 的作用

## C.2 CET 与教育技术的对话

与脚手架理论（**Scaffolding Theory Wood Bruner & Ross 1976**）的关系 核心主张：- 教学应提供临时性支持（脚手架）- 随着学习者能力提升逐步撤出支持 - 最终学习者完全独立

CET 的立场：

相似点（我们如何借鉴）：- 完全一致：EML 的系统性支持削减正是脚手架理论的形式化 - 借鉴了 Vygotsky 的最近发展区（ZPD）概念（50-70% 挑战区）

差异点（我们如何超越/补充）：- CET 的贡献：1. 量化削减过程：将”渐进撤出形式化为削减曲线  $S”(t)$  2. 提供验证标准：AVP 判据确保撤出成功（不是撤了就行要看  $P_2$ ）3. 扩展到 AI 时代：自动化脚手架（LSA 架构见第 5 章）4. 跨尺度扩展：从个体学习扩展到团队/组织能力建构

互补点（如何协同）：- 脚手架理论提供教育学基础 - CET 提供 AI 时代的技术实现路径 - 未来合作方向：- AI 如何自动判断学习者准备好了（能力监测  $C(t)$ ）- 不同学科的最优脚手架设计（跨领域校准）

对话焦点：

脚手架理论：何时撤出支持 CET 回应：基于能力监测（ $C(t)$ ）动态决定并通过 AVP 验证撤出是否成功

脚手架理论：如何知道学习者准备好了 CET 回应：监测成功率（目标 50-70% 区间工作假设需跨领域/任务校准）和能力向量  $C(t)$  的变化趋势

脚手架理论：如果撤出失败怎么办 CET 回应：自动回退机制（见 5.3.3 节 SGS 的安全约束）如连续失败触发红色预警  $S(t)$  回升至安全水平

## C.3 CET 与 AI 伦理的对话

与自主性（**Autonomy**）讨论的关系 AI 伦理的核心关切：- AI 是否威胁人类自主性 - 如何保护用户的选择自由 - 算法决策的透明性和可解释性

CET 的立场：

相似点（我们如何借鉴）：- 共同关注”人类主体性”（human agency）- 认同用户应保持对 AI 的控制

差异点（我们如何超越/补充）：- CET 超越传统自主性讨论：不仅关注”能否选择”更关注能力是否保持 - 认知自主性（Cognitive Autonomy）：- 传统自主性：我可以选择用或不用 AI（意志自由）- 认知自主性：我用了 AI 后仍然有能力独立完成任务（能力自由）

CET 的独特贡献：1. 将抽象的自主性转化为可测量的独立能力（AVP）2. 提出伙伴式主体性作为 AI 角色的伦理基准（见 3.0.4 节）3. 警示外骨骼模式对认知自主性的长期威胁

互补点（如何协同）：- AI 伦理提供规范性框架（应该如何）- CET 提供可操作的评估工具（如何验证）- 未来合作方向：- 将 AVP 纳入 AI 系统的伦理审查流程 - 开发认知自主性指数”作为 AI 产品的标签”（类似能效标签）

对话焦点：

AI 伦理：AI 威胁人类自主性吗 CET 回应：取决于 AI 的设计外骨骼模式威胁认知自主性（ $P_2 < B_0$ ）内共生模式增强认知自主性（ $P_2 > B_0$ ）

AI 伦理：如何量化自主性 CET 回应：通过 AVP 判据如果  $P_2 > B_0$  说明用户的认知自主性得到维护甚至增强

AI 伦理：用户知情同意就够了吗 CET 回应：不够即使用户同意使用 AI 如果导致能力退化（外骨骼）从长期伦理看是有问题的需要能力保护机制（EML 设计）

#### C.4 CET 与组织行为学的对话

与组织韧性（**Organizational Resilience Hollnagel 2011**）的关系 核心主张：- 韧性 = 在扰动下维持功能的能力 - 需要多样性、冗余、适应性 - 复杂系统的”韧性工程”（Resilience Engineering）

CET 的立场：

相似点（我们如何借鉴）：- 完全一致：O-AVP 正是组织认知韧性的测量 - 借鉴了多样性（角色冗余）和冗余（知识流动）概念

差异点（我们如何超越/补充）：- CET 识别了 **AI** 依赖作为韧性的新威胁：- 传统威胁：单点故障、人员流失、供应链中断 - 新威胁：AI 依赖导致的集体能力退化 - CET 提供了具体测量方法：- 48h 演练（见 4.2.2 节）- BCI/ICR 双指标（见 4.2.3 节）

互补点（如何协同）：- 组织韧性理论提供宏观框架 - CET 提供 AI 时代的具体威胁识别和测量工具 - 未来合作方向：- 将 O-AVP 纳入组织风险评估体系 - 开发认知韧性仪表盘（real-time monitoring）

对话焦点：

组织韧性理论：如何增强组织韧性 CET 回应：除了传统措施（多样性、冗余）还需关注”认知韧性”——确保组织在 AI 不可用时仍能运作

组织韧性理论：如何测量韧性 CET 回应：O-AVP 提供了可操作的测量方法：48h 中断演练 + 恢复时间/独立完成率

组织韧性理论：韧性是否需要持续投入 CET 回应：是的建议每季度一次 O-AVP 演练（类似消防演习）保持组织的肌肉记忆

#### C.5 CET 与系统科学的对话

与复杂系统理论（**Complex Systems Theory**）的关系 核心主张：- 复杂系统表现出涌现性（emergence）、非线性、自组织 - 局部优化 全局最优 - 反馈环路主导系统行为

CET 的立场：

相似点（我们如何借鉴）：- 认同跨尺度涌现性：I-AVP 通过 T-AVP 必然通过（见 4.1 节）- 借鉴反馈环路思想：AVP 闭环（见 5.4 节）

差异点（我们如何超越/补充）：- CET 关注认知系统的特殊性：- 能力的不可逆性（退化容易恢复难）- 代际传承的路径依赖（ $T_0 \rightarrow T_1 \rightarrow T_2$ ）- 认知公地悲剧（个体理性  $\rightarrow$  集体退化见 4.3 节）

互补点（如何协同）：- 复杂系统理论提供分析框架 - CET 提供 AI-人类系统的具体案例 - 未来合作方向：- 建立 CET 的系统动力学模型（SD modeling）- 模拟长期演化路径（agent-based modeling）

对话焦点：

复杂系统理论：为何局部优化（个体用 AI 效率高）导致全局问题（组织能力退化）CET 回应：因为存在认知公地悲剧每个人理性选择使用 AI（ $P_1$  高）但集体结果是能力退化（ $P_2 < B_0$ ）类似过度捕捞的公地悲剧

复杂系统理论：如何避免系统崩溃 CET 回应：通过多尺度监测（I/T/O/S-AVP）和早期预警（黄色/红色预警见 5.4.4 节）在系统接近临界点前干预

复杂系统理论：系统是否有”吸引子”（attractor）CET 回应：可能存在两个吸引子：- 良性吸引子：内共生平衡态（ $P_2 \geq B_0$  持续稳定）- 病理吸引子：外骨骼锁定态（ $P_2 < B_0$  且难以恢复）需要实证研究验证这一假设

附录 D：术语与符号索引

D.1 核心概念术语（按字母序）

认知内共生（**Cognitive Endosymbiosis**）- 首次出现：1.3.1 节 3.0.4 节（完整定义）- 核心定义：AI 作为认知伙伴通过有益摩擦和支持削减促进用户独立能力提升的人机协作模式（见 3.0.3 节 EML 定义锚点 B2）- 关键特征：有益认知摩擦 系统性支持削减 AVP 验证 - 对照术语：认知外骨骼（病理模式）- 参见：3.3 节（EML）、4.1 节（跨尺度扩展）

认知外骨骼（**Cognitive Exoskeleton**）- 首次出现：1.2.2 节（问题提出）3.0.6 节（详细定义）- 核心定义：过度依赖 AI 导致独立能力退化的病理模式特征是高  $P_1$ （有 AI 表现好）但低  $P_2$ （无 AI 表现差）- 关键特征：零摩擦设计 无支持削减 AVP 失败（ $P_2 < B_0$ ）- 对照术语：认知内共生（目标模式）- 参见：表 3.5（失败模式分类）、附录 B.1（案例）

反脆弱性验证原则（**Antifragility Validation Principle AVP**）- 首次出现：1.3.1 节（概念引入）3.0.2 节（完整定义）- 核心定义：以拔线测试检验协作是否促进独立能力判据为  $P_2 \geq B_0$ （见 3.0.2 节 AVP 定义锚点 B1）- 关键参数： $B_0$ （基线）、 $P_2$ （拔线后能力）、（最小提升阈值）、W（拔线窗口）- 分级体系：AVP-Basic（ $P_2 \geq B_0$ ）、AVP-Retention（ $P_2 \geq B_0$ ）、AVP-Transfer（+ 迁移）- 参见：3.1 节（详细机制）、附录 A（测量工具包）

内共生最小法则（**Endosymbiotic Minimal Law EML**）- 首次出现：1.3.2 节 3.0.3 节（完整定义）- 核心定义：构成认知内共生的设计必要条件：有益认知摩擦（50-70% 成功率） 系统性支持削减（ $S_4 \rightarrow S_1 \rightarrow S_0$ ）（见 3.0.3 节 EML 定义锚点 B2）- 与 AVP 关系：EML 是设计条件 AVP 是验收条件二者联合构成充分必要条件 - 参见：3.2 节（有益摩擦）、3.3 节（支持削减）

伙伴式主体性（**Partner-like Agency**）- 首次出现：1.3.3 节 3.4 节（详细阐述）- 核心定义：AI 作为认知伙伴的理想角色定位三个锚点：摩擦注入 脚手架消退 AVP 闭环 - 关键特征：功能性非拟人化教练模式而非仆人模式 - 参见：3.4 节（理论阐述）、5.5 节（伦理治理）

有益认知摩擦（**Beneficial Cognitive Friction**）- 首次出现：3.2 节 - 核心定义：适度挑战（目标成功率 50-70% 工作假设需跨领域/任务校准）促进能力增长的设计策略 - 类型：完整性摩擦、抽象性摩擦、延迟性摩擦 - 理论基础：心流理论（Csikszentmihalyi）+ 最近发展区（Vygotsky）- 参见：表 3.4（摩擦类型对照）、5.2 节（CFE 引擎）

系统性支持削减（**Systematic Support Reduction**）- 首次出现：3.3 节 - 核心定义：AI 支持强度按既定削减曲线从  $S_4 \rightarrow S_1 \rightarrow S_0$  的设计策略 - 削减曲线类型：线性、指数、S 型 - 安全机制： $S_{\min}$  下限、回退触发、预警系统 - 参见：图 B.3（削减曲线对比）、5.3 节（SGS 调度器）

拔线测试（**Unplugged Test**）- 首次出现：3.1 节 - 核心定义：在无 AI 环境下测量独立能力的标准化测试方法 - 关键要素：等值性（任务难度相同）、拔线窗口 W、违规监测 - 参见：附录 A.1.2（测量执行协议）

D.2 符号与参数索引

能力与表现：

符号	含义	单位/范围	首次定义
$B_0$	基线能力（使用 AI 前）	任务特定评分	3.0.2 节
$P_1$	协作期表现（使用 AI 时）	任务特定评分	3.0.2 节
$P_2$	拔线后独立能力	任务特定评分	3.0.2 节
$C(t)$	能力向量（多维）	n 维向量	5.4.3 节
$\Delta C$	能力增长	$P_2 - B_0$	3.0.2 节

阈值与参数:

符号	含义	默认值	首次定义
	最小有意义提升阈值	Cohens d 0.3 或 10% ( <i>working assumption</i> ) *	3.0.2 节
W	拔线窗口	4-8 周 (默认 6 周) *	3.0.2 节
F	摩擦参数	0-1 (目标 0.5-0.7) *	3.2 节
S(t)	支持强度函数	0-1 从 S0→0*	3.3 节
S0	初始支持强度	0.8 (80%) *	3.3.1 节
S_min	安全支持下限	0.2 (20%) *	3.3.2 节
	削减速率	任务特定 *	3.3.1 节

\* (工作假设需跨领域/任务校准)

尺度特定指标:

符号	含义	阈值	首次定义
I-AVP	个体反脆弱性验证	$P_2 B_0+$	4.1 节
T-AVP	团队反脆弱性验证	0.7 (群体级) *	4.1.3 节
O-AVP	组织反脆弱性验证	告警 0.70 目标 0.85*	4.2.3 节
BCI	业务连续性指数	0-1	4.2.3 节
ICR	独立完成率	0-1	4.2.3 节
S-AVP	社会层认知资本指标	代际对比	4.3 节

LSA 架构:

符号	含义	功能	首次定义
L1	基础 AI 能力层	AI 模型接入、知识库	5.1 节
L2	摩擦与削减层	CFE+SGS	5.2-5.3 节
L3	监测与反馈层	AVP-TM+ 预警系统	5.4 节
L4	编排与治理层	多尺度协调 + 伦理治理	5.5 节
S4→S1→S0	支持档位栈	强 → 弱的 4 级支持强度	3.3 节
CFE	认知摩擦引擎	动态摩擦调节	5.2 节
SGS	支持削减调度器	削减曲线管理 + 回退	5.3 节
AVP-TM	AVP 遥测模块	能力监测 + 数据采集	5.4 节
MSO	多尺度编排器	I/T/O 跨尺度协调	5.5 节

D.3 缩写速查表 (完整版)

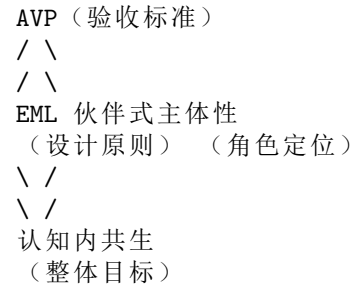
缩写	全称	中文	首次出现
AVP	Antifragility Validation Principle	反脆弱性验证原则	1.3 节
EML	Endosymbiotic Minimal Law	内共生最小法则	1.3 节
LSA	Layered Symbiosis Architecture	分层共生架构	1.3 节
LSA-F	LSA Functional Hierarchy	LSA 功能分层	3.0.5 节
CFE	Cognitive Friction Engine	认知摩擦引擎	5.2 节
SGS	Support Graduation Scheduler	支持削减调度器	5.3 节



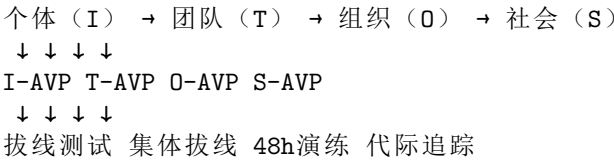
缩写	全称	中文	首次出现
AVP-TM	AVP Telemetry Module	AVP 遥测模块	5.4 节
MSO	Multi-Scale Orchestrator	多尺度编排器	5.5 节
I-AVP	Individual AVP	个体反脆弱性验证	4.1 节
T-AVP	Team AVP	团队反脆弱性验证	4.1 节
O-AVP	Organizational AVP	组织反脆弱性验证	4.2 节
S-AVP	Societal-level indicators	社会层认知资本指标	4.3 节
BCI	Business Continuity Index	业务连续性指数	4.2 节
ICR	Independent Completion Rate	独立完成率	4.2 节
IRT	Item Response Theory	项目反应理论	附录 A
RCT	Randomized Controlled Trial	随机对照试验	6.2 节
DID	Difference-in-Differences	差分中的差分	6.2 节
ZPD	Zone of Proximal Development	最近发展区	2.3 节
WEIRD	Western Educated Industrialized Rich Democratic	西方、受教育的、工业化的、富裕的、民主的	6.1 节

#### D.4 跨章节概念地图

核心三角关系：



跨尺度扩展链：



LSA 技术实现栈：



参考文献