

# 神经符号人工智能的分层协同架构 及其跨领域工程实践印证

## Hierarchical Collaborative Architecture for Neuro-Symbolic AI and Cross-Domain Engineering Validation

2026年1月3日凌晨，委内瑞拉首都加拉加斯，一场闪电般的军事行动震惊世界。美军成功突袭并抓捕了总统尼古拉斯·马杜罗，整个过程行云流水，几乎没有给对手留下任何反应时间。事后披露的信息揭示了一个关键事实：这次“斩首”行动的成功，并不仅仅是特种作战的胜利，更是一次由人工智能驱动的“算法定义战争”的完美演示。隐藏在背后的，是 Palantir 为美军量身定制的 Maven Smart System——一个能够整合海量战场数据、进行实时推理与决策的智能中枢。它让美军指挥官得以穿透迷雾，看见一个由符号化实体和关系构成的“数字战场”，并迅速完成从情报到打击的闭环。



图 1 美军突袭并抓捕委内瑞拉总统尼古拉斯·马杜罗。

这一事件向世界发出强烈信号：在真实、复杂、高风险的物理世界中，人工智能系统必须具备可解释、可控制、可验证的能力。而这，正是神经符号人工智能（NSAI）的核心追求——通过“感知—认知—行动”的分层协同架构，在强大的神经感知层与精准的控制执行层之间，插入一个**显式的、符号化的认知与约束中间层**，系统性地赋予智能系统可解释、可控制、可验证的关键属性。这一理念已超越理论探讨，在自动驾驶、绿色制造、企业软件、战场空间等四大差异显著的复杂领域，获得了详尽而坚实的工程验证，共同指向构建下一代可信赖智能的统一范式。

## 一、神经符号人工智能的分层协同架构

神经符号人工智能的设计哲学源于一个根本洞察：纯粹的神经网络（如大模型）擅长从数据中学习模式，但其决策过程如同黑箱，难以在关键任务中被信赖；而纯粹的符号系统虽逻辑透明，却难以处理非结构化的现实世界。NSAI 将二者深度融合，构建了一个“感知—认知—行动”的三层递进架构：

1. **感知层（神经）**：负责处理多模态原始输入（图像、语音、传感器数据等），利用深度神经网络提取高层特征，并将其转化为初步的语义表示。这一层的目标是实现“世界是什么样”的高效、鲁棒感知。

2. **认知层（符号）**：这是 NSAI 架构的核心。它接收感知层输出的语义信息，并结合领域先验知识（如物理定律、业务规则、交通法规），通过显式的符号化推理生成可解释的决策依据。具体而言，该层包含：

- **符号化状态表征**：将连续、高维的感知数据映射为离散、可操作的符号断言（例如“前方车辆减速”、“螺钉已对准”、“库存不足”）。

- **知识与规则引擎**：内置形式化的知识库（本体、规则、约束），用于对符号状态进行逻辑推理、规划与验证。

- **决策可追溯机制**：所有推理步骤均被记录，使得决策过程透明、可审计。

3. **行动层（控制）**：将认知层输出的符号化指令转化为精确的、符合物理约束的控制信号，驱动执行器完成具体动作。该层常内嵌经典控制理论模型（如 PID、阻抗控制），确保操作的稳定性与安全性。

这一架构的核心价值在于：**神经系统赋予机器适应开放世界的灵活性，符号系统则为其植入可被理解、可被信任的“骨架”与“护栏”**。当面对未知场景时，系统不是盲目猜测，而是基于符号化推理和先验知识进行稳健决策，从而从根本上破解了“黑箱”困境。

## 二、自动驾驶：英伟达 AlpaMayo——开放环境中的可解释决策



图 2 英伟达推出 AlpaMayo 系列开源 AI 模型与工具，加速安全可靠的推理型辅助驾驶汽车开发。

### 1. 关键问题：黑箱决策与可观察性塌缩

传统端到端自动驾驶模型将感知、决策与控制压缩至单一神经网络，导致两个根本性缺陷：其一，决策过程如同“黑箱”，系统无法回答“为何如此决策”的质询，难以进行安全认证与责任界定；其二，模型潜在表示存在“可观察性塌缩”，倾向于保留视觉显著特征，而忽略或混淆对控制至关重要的物理量（如他车意图、风险概率），导致在未知场景中基于失真状态决策，可靠性存疑。

## 2. NSAI 思想：引入显式推理层的分层架构

英伟达于 2026 年消费类电子产品展（CES）发布并开源的自动驾驶模型 Alpamayo，从工程实践的角度直接回应了上述问题。其核心是设计了一个清晰的“视觉—推理—行动”三层架构：

- **视觉层（神经感知）**：利用 Transformer 等模型处理多模态传感数据，完成高效的场景理解与特征提取。
- **推理层（符号核心）**：作为一个独立的“推理骨干网络”，接收结构化感知信息，并结合交通规则、导航目标进行**因果推理、意图预测与逻辑判断**。其输出是**可解释的“思维链”**，能将决策依据（如“因前车减速，故保持车距”）显式化、结构化。
- **行动层（受控执行）**：基于推理层输出的符号化目标与约束，生成具体、平滑且符合车辆动力学的控制轨迹。

## 3. 工程印证价值：实现从“验证结果”到“验证逻辑”的范式转变

Alpamayo 通过引入**显式、可追溯的推理层**，成功将驾驶智能从“黑箱”推向“白箱”。它印证了在安全攸关的开放动态环境中，**可解释性是可信赖的前提**。配合其开源仿真框架 AlpaSim，开发者得以对决策逻辑进行反事实测试与多策略回放，标志着自动驾驶系统的验证标准，从过去只关注“行为结果是否正确”，

升级为同时检验“决策逻辑是否合理”，为高阶自动驾驶的安全部署奠定了关键的工程基础。

### 三、绿色制造：动力电池自主拆解——非结构化场景中的可信智能体



图 3 动力电池多机器人自主协作拆卸系统 BEAM-1

#### 1. 关键问题：动态不确定性下的安全与可靠性危机

退役动力电池拆解是典型的“非结构化”工业任务，面临三重不确定性：对象不确定性（型号繁多、磨损各异）、环境不确定性（螺钉锈蚀、连接件变形）、任务不确定性（需实时规划拆卸顺序与力度）。传统预编程或纯学习驱动的机器人，在此类场景中缺乏对关键物理状态的理解与安全边界的认知，极易引发碰撞、损坏甚至安全事故，可靠性无法保障。

#### 2. NSAI 思想：构建“符号感知—原语规划—闭环控制”的融合系统

社区围绕 BEAM-1 机器人等平台，构建了完整的 NSAI 解决方案：

- 符号化状态感知：通过“神经谓词”模型，将原始视觉信息实时编译为如“螺钉-已定位-未对准”、“连接件-处于卡紧状态”等符号断言。这相当于为机

器人生成一份可读的、关于当前工作场景的“事实清单”，从根本上提升了系统对隐藏关键状态的可观察性。

- **知识与原语驱动的规划**: 系统包括一个融合拆卸工艺知识的参数化动作原语库（如“力控套接”、“柔顺分离”）。符号规划器基于当前状态事实，通过逻辑推理调用并序列化原语。每个原动作语都封装了明确的前置条件、后效状态及物理安全约束，严格定义了机器人的可控制行为边界。

- **感知—控制闭环与先验嵌入**: 在执行层面，基于阻抗控制的柔顺控制模型被显式嵌入，使机器人能“感知”力并动态适应物理交互中的不确定性，保障操作的稳定性。同时，高保真数字孪生系统提供了“仿真—验证—优化”的闭环，使安全成为可提前测试的结构属性。

### 3. 工程印证价值：验证 NSAI 在高危工业场景中的可靠性工程路径

该实践成功证明，NSAI 架构能将神经网络处理非结构化感知的灵活性，与符号系统施加精确物理约束的可靠性完美结合。它使工业机器人从被动执行固定程序的工具，转变为能主动理解环境、解释自身计划并确保行为安全的智能体。这为在绿色制造、危化品处理等需要极高可靠性的领域实现“自主化”，提供了经得起验证的技术路径与工程范本。

## 四、企业软件：Palantir AIP——复杂商业系统中的本体化约束

### 1. 关键问题：商业 Agent 的“失控”与决策失真

在企业运营中，直接基于大模型构建的智能体（Agent）常陷入“能用但不敢用”的困境。其核心问题是：Agent 直接操作杂乱的数据库、API 和文档，缺乏对业务实体（如“库存”、“订单”）及其复杂关系、规则的一致理解。这导

致其决策语义失准、行为越界、逻辑不可追溯，无法承担核心业务流程的自动化职责，信任难以建立。

## 2. NSAI 思想：以“企业级本体”构建符号化业务现实

Palantir AIP 的核心理念是引入“本体 (Ontology)”作为企业数据的符号化中间层。

- **统一语义，构建可观察的业务镜像**：Ontology 形式化地定义企业内所有人、物、事件、流程及其关系（如“可用库存=总库存-锁定库存”）。它为 Agent 提供了一个无歧义的、符号化的业务世界模型，彻底解决了因数据孤岛和语义混淆导致的理解偏差。

- **定义受约束的行为空间**：所有业务操作（如“审批采购单”、“调拨库存”）都在本体中被定义为受权限和业务规则严格约束的**动作**。Agent 只能在本体规定的“合法动作空间”内行事，其**可控制性**得到了根本保障。

- **内生可追溯性**：任何决策都能自动关联到其所依据的原始数据、触发的业务规则及影响的对象，实现**全链路审计追踪**，决策过程从黑箱变为白箱。

## 3. 工程印证价值：证明 NSAI 原则在信息空间的普适性与商业必然性

Palantir 在金融、供应链、政府等高风险领域的成功，雄辩地印证了 NSAI 思想的普适性。它表明，即使在没有物理实体的纯信息领域，智能体的“可信”问题根源同样在于缺乏一个符号化的、富含知识的**认知中间层**。通过构建企业本体，AIP 实现了商业决策的**可观察、可控制、可验证**，使得 AI Agent 得以深度融入核心业务系统。这证明，NSAI 不仅是自动驾驶、机器人技术的选项，更是构建未来**可信赖企业级智能的工程必然**。

## 五、军事领域：Palantir Maven Smart System —— AI 时代战争的“决策中枢”



图 4 Palantir Maven 智能系统将 AI 算法定义战争从概念推向前线。

2026 年 1 月 3 日凌晨，美军入侵委内瑞拉并抓捕总统马杜罗，引发国际社会强烈谴责。在委内瑞拉总统马杜罗突袭事件震惊世界的背后，是 NSAI 分层协同架构在军事领域的巅峰实践。Palantir 为美军量身定制的 **Maven Smart System (MSS)**，现已成为美军实现“联合全域指挥与控制”（CJADC2）战略的基石，将“**算法定义战争**”从概念推向前线。

### 1. 关键问题：战场数据的“混沌”与决策延迟

现代战场数据呈现海量、多源、异构的典型特征：卫星图像、信号情报、线人报告、社交媒体信息每秒都在涌入，但这些数据彼此孤立、格式混乱、语义不明。传统情报分析依赖人工，速度慢且易出错，导致从“传感器到士兵”的链路存在致命延迟。其根源在于，系统缺乏对战场实体（如“马杜罗”、“军事堡垒”、“黑色交通工具”）及其复杂关系的统一、可计算模型，使得决策如同在迷雾中导航，无法实现实时、精准的打击闭环。

## 2. NSAI 思想的映射：以“Ontology+知识图谱”构建数字战场

Palantir MSS 的核心理念，正是 NSAI 控制论立场在军事领域的完美映射。

其通过构建企业级本体与知识图谱，系统性地解决了战场认知难题：

- **Ontology: 为战场数据定义“通用语言”**。它将所有杂乱原始数据，强制转化为结构化的“实体”类型，如“人物（马杜罗）”、“地点（军事堡垒）”、“事件（通信中断）”。这相当于为战场建立了一本统一的“数据字典”，从根本上提升了系统对关键作战要素的可观察性，使 AI 不再“看见”数据噪声，而是“观察”到明确的符号化事实。

- **知识图谱：编织战场“关系网络”**。一旦数据被结构化，知识图谱便建立实体间的动态关系，如“马杜罗‘位于’军事堡垒”、“CIA 线人‘报告了’马杜罗的行踪”。这使得 AI 能够进行多跳推理与模式识别，自动将碎片化信息关联到唯一的“数字马杜罗”，为指挥官呈现一个实时更新、三维立体的可控制、可追溯的战场镜像。

- **AIP 大模型引擎：赋能战略决策**。在 Ontology 与知识图谱之上，Palantir 的 AIP（人工智能平台）集成了大模型能力，扮演“超级情报分析师”与“虚拟战争推演室”的双重角色。它能极速阅读并摘要万份情报，也能模拟数百万次突袭路径，生成附带风险评估的最优方案，并自动进行交战规则 (ROE) 合规检查，将决策过程从数周的“人工会议”压缩为实时的“人机对话”，确保了军事行动在复杂约束下的可控制性与可验证性。

## 3. 工程印证价值：NSAI 原则在军事领域的决定性胜利

马杜罗突袭行动的成功，是 NSAI 架构在最高风险、最复杂环境中获得的“实战认证”。它雄辩地证明：无论是操控无人机群（协同作战飞机 CCA）、

实现毫秒级目标识别 (ATR) ，还是规划跨军种联合突袭，构建可信作战智能的关键路径是一致的——必须通过一个符号化的、富含知识的认知中间层 (Ontology/知识图谱) 来桥接原始感知 (情报数据) 与最终行动 (打击指令) ，从而系统性地保障战场态势的可观察性、作战决策的可控制性与打击闭环的稳定性。

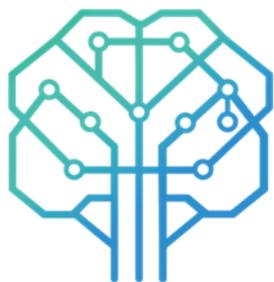
这一实践标志着，NSAI 不仅是工业与商业智能的选项，更是决定未来战争胜负的国家战略级基础设施。它将 AI 时代的军事斗争推向了“算法定义战争”的新纪元。

### **结束语：迈向统一的可信赖智能新范式**

从开放道路 (Alpamayo)、工厂车间 (动力电池拆解)、全球供应链 (Palantir AIP) 到战场空间 (Maven Smart System) 四个案例虽领域迥异，却共同验证了同一套神经符号分层协同架构的强大生命力。它们清晰地表明，无论是应对物理不确定性、商业复杂性还是战场对抗性，在数据驱动的“感知/行动”闭环中，嵌入一个承担状态符号化、逻辑推理与边界定义职能的认知层，是破解人工智能系统“不可解释、不可控制、不可验证”困境的根本出路。

这标志着人工智能的发展正进入一个新时期：从一味追求“更大规模、更宽泛化”的模型能力竞赛，转向追求“更高可信度、更强责任感”的系统工程深化。神经符号人工智能所提供的分层协同架构，正是引领这一范式变革、构建能在真实世界中担当重任的可信赖智能的统一蓝图。

**神经符号AI，赋能绿色制造的人工智能引擎**  
**<https://www.nsaihome.org.cn>**



**NSAiHome**  
神经符号人工智能社区