

# Agent 隐蔽通信 内容与行为的分层协同与动态双信道模型

参考url : <https://github.com/precize/OWASP-Agentic-AI/blob/main/agent-covert-channel-exploitation-16.md>

cursor\_name : Covert Channel Exploitation in AI Agent

## ▼ 我目前研究方向和 idea

### ▼ 现有 LLM 内容隐写的“静态”困境

- **问题核心：同步悖论。** 现有的基于 LLM 的隐写方法，尤其是依赖特定 Prompt 模板的技术，在需要更新或优化 Prompt 时面临两难：
  - 直接用旧 Prompt 生成的内容来传递新 Prompt，会暴露通信意图（因为内容可能效果不佳或与新 Prompt 不匹配）。
  - 不传递新 Prompt，则双方无法同步，隐写通信无法继续或效率低下。
- **限制：** 这种对 Prompt 的强依赖性导致了同步难题，限制了 LLM 隐写的灵活性和鲁棒性。

### ▼ 提出 Agent 的潜力未被挖掘

- LLM 作为 Agent 使用时，其能力远不止于生成文本。
- Agent 可以与环境交互、感知状态、进行推理，并根据目标调整自身行为。

**关键洞察：** Agent 的行为本身可以作为除内容之外的第二个潜在信道。

### ▼ 明确核心挑战

- 如何设计一个系统，能够同时利用 Agent 的内容生成能力（高带宽，潜在高隐蔽性）和行为交互能力（用于同步和控制），构建一个协同工作的隐写机制？
- 目标是克服静态 Prompt 的同步瓶颈，最终实现 Agent 在复杂环境中自主、动态、隐蔽地进行信息传递。

### ▼ 提出核心理念：内容与行为的分层协同与动态双信道模型

- **双信道模型：**
  - **内容信道：** 利用 LLM 生成的文本作为载体，负责传输主要的、**大容量的**秘密信息。这是传统隐写的主要关注点。
  - **行为信道：** 利用 Agent 在模拟环境中的**交互行为模式**（如消息发布频率、时间间隔、点赞/评论特定内容、互动对象选择等）作为载体。这个信道负责传输**低容量、但关键的控制/同步信号**。
- **解决同步悖论的关键：** 行为信道用于传递关于内容信道状态的元信息，特别是 **Prompt 的变更信号**。
  - 例如，通过一个简短的行为信号 (如 `Prompt1 → Δprompt → 20bits → Prompt2` 的示意)，来指示内容信道需要切换到新的 Prompt。

## ▼ 设计具体的协同机制（以行为信道传递 Prompt 变更信号为例）

- **信号灯机制 (发送方 Alice)：**
  - 当 Alice 需要更换 Prompt 时，她不会直接将新 Prompt 嵌入到内容信道（因为这可能不安全或暴露意图）。
  - Alice 使用新的 Prompt 生成内容文本（发布在内容信道）。
  - **同时**，Alice 控制其 Agent 在行为信道上展现出一种**特殊、预定义的行为模式**，这个行为模式本身就编码了“Prompt 已变更”的信号。
- **接收方 (Bob) 解码逻辑：**
  - Bob 的 Agent 首先监控 Alice Agent 的行为。
  - 一旦检测到预设的“Prompt 变更”行为信号，Bob 就知道接下来收到的内容文本需要使用新的（或某种预定义的备用/更新）Prompt 来解码。
  - 如果未检测到行为信号，则继续使用旧 Prompt 解码。

## ▼ Δ prompt

### ▼ 行为信号指示 Prompt 切换到底在传递什么

1. 例如，假设系统最初使用的是一个代号为 `初始指令集` 的 Prompt。
2. 当发送方 Agent 需要让接收方 Agent 切换到一个新的 `后续指令集` 时，它不会直接在常规的内容信道发送“请切换到后续指令集”这样的明文。
3. 相反，发送方 Agent 会执行一个预先约定好的、特定的行为模式。

4. 比如，这个行为可以是：‘在接下来的连续三次常规状态更新中，Agent 发送更新消息的时间戳尾数（例如毫秒数）遵循一个特定的、预先约定的非随机模式（例如，依次为奇数、偶数、奇数）。与此同时，这三次更新中，第二次更新的内容长度相比平时平均长度有一个微小但统计上显著的增加（例如，增加5-10个字符，内容本身依然是合法的常规更新）’
5. 接收方的 Agent 持续监测发送方 Agent 的这些可观察行为。
6. 一旦它识别出这个预设的、略显反常的行为组合（这个行为组合本身可以被编码为一个极短的信号），接收方就知道，接下来通过内容信道发送过来的信息，应该使用双方已知的 **后续指令集**（或者根据信号指示从某个预定义的列表中选择新的指令集，或者准备接收一个通过内容信道隐蔽传输的 **后续指令集** 的更新包）来进行解析和执行。

行为信道传递的是一个**高层级的控制指令或同步信号**（例如‘切换到编号为X的Prompt’或‘启用紧急模式Prompt’），而不是 **后续指令集** 的全部文本内容。这个控制信号因为信息量小，所以非常适合通过低带宽、难以察觉的行为信道进行传递。”

- **信号是具体的行为模式**：而不仅仅是抽象的数字或比特。
- **信号的含义**：是控制指令，用于同步双方对内容信道所用 Prompt 的认知。
- **信道特性**：强调了行为信道低带宽、高隐蔽性的特点，适合传输短小的控制信号。

## ▼ 如何实现 prompt 切换

### ▼ 预设 Prompt 组，然后选择/变化：

通过预先定义好若干组 Prompt 或 Prompt 模板，可以将“完全传输新 Prompt”的问题转化为“选择哪个 Prompt/模板”或“如何微调已知 Prompt”。

这大大降低了需要实时传递的信息量。

### ▼ 如果使用 VAE (Variational Autoencoder) 压缩潜空间

- **想法**：将 Prompt 编码到 VAE 的潜空间，然后只传递潜向量的变化 ( $\Delta z = z_{new} - z_{old}$  或者直接传递  $z_{new}$ )。
- **优点**：
  - 如果潜空间维度远小于原始 Prompt 的表示维度，可以实现高压缩。

- 潜空间可能捕获 Prompt 的语义信息，理论上小的潜向量变化可能对应有意义的 Prompt 语义变化。

- **挑战与思考：**

- **有损压缩与精度：**VAE 通常是有损的。对于 Prompt 而言，即使是微小的重建误差（比如一个词的改变、语气的细微偏差）也可能导致 Agent 行为的巨大差异甚至完全错误。
- 需要评估这种有损性是否可以接受。如果目标是精确同步 `prompt2`，VAE 可能不是最佳首选，除非重建质量非常高。
- **模型共享与预训练：**发送方和接收方都需要拥有完全相同的、预训练好的 VAE 模型。训练一个能很好泛化各种 Prompt 的高质量 VAE 本身就是一项挑战。
- **潜向量变化的大小：**即使是潜向量的变化  $\Delta z$ ，其信息量（需要多少比特来表示）是否足够小以适应行为信道（提到“行为空间可以嵌入的量应该不大吧”）？这取决于潜空间的维度和表示精度。
- **线性假设：**简单地传递  $z_{new} - z_{old}$  依赖于潜空间中 Prompt 语义变化的某种线性或简单可加性，这不一定成立。

## ▼ 除了 VAE，还有哪些方法可以考虑？

### 1. 索引化/ID 化 Prompt + 差分编码 (针对组内微调)：

- **工作方式：**

- 双方共享一个预定义的 Prompt 库（或者“Prompt 组”）。每个 Prompt 或 Prompt 模板都有一个唯一的 ID 或索引。
- **选择新的基础 Prompt：**如果 `prompt2` 是库中一个全新的 Prompt，只需通过行为信道传递其 ID。这是最高效的方式。
- **微调现有 Prompt：**如果 `prompt2` 是基于库中某个已知 Prompt (`prompt_base`，通过 ID 选定) 进行的修改，则计算  $\Delta \text{prompt} = \text{diff}(\text{prompt\_base}, \text{prompt2})$ 。这个 `diff` 可以是文本差异（如 `git diff` 的输出格式），也可以是更结构化的差异描述。

- **压缩  $\Delta \text{prompt}$ ：**

- **稀疏性利用：**如果  $\Delta \text{prompt}$  只是少数几处小的修改，可以设计非常紧凑的格式来描述这些修改（例如：行号、旧词、新词；或特定参数的变更）。

- **通用无损压缩**：对生成的 `diff` 文本应用如 Huffman 编码、LZ 算法等进一步压缩。
- **预定义编辑操作码**：如果修改类型有限（例如：替换关键词、增删特定指令、调整某个权重参数），可以为这些操作预定义极短的操作码。

## 2. 参数化 Prompt 模板：

- **工作方式**：如果“Prompt 组”实际上是带有可变参数的 Prompt 模板（例如：“请以{风格}风格，就{主题}生成一段{长度}的文本”）。
- **传输**：当从 `prompt1` 切换到 `prompt2` 时，如果只是模板中的参数值改变了，行为信道仅需传递这些改变了的参数的新值。
- **压缩**：参数值（特别是如果是枚举型、布尔型或小范围整数）本身就很容易压缩，甚至可以用几比特表示。

## 3. 语义级别的差分与编码：

- **工作方式**：不同于纯文本 `diff`，这种方法尝试理解 Prompt 的语义结构。

例如，将 Prompt 解析成一种内部表示（如意意图、指令序列）。`Δprompt` 就是这种内部结构上的变化。

- **传输**：传输描述这些结构性变化的指令。
- **优点**：可能比文本 `diff` 更紧凑，更能抓住核心变化。
- **挑战**：需要更复杂的解析和生成逻辑，双方对语义结构需有一致理解。

## 4. 基于上下文的增量更新与预设“补丁”：

- **工作方式**：针对常见的 Prompt 演变模式，预先准备好一些“补丁集”或“转换规则集”。
- **传输**：行为信道传递一个信号，指示接收方应用哪个预设的补丁集或转换规则到当前的 `prompt1` 来生成 `prompt2`。
- **例子**：信号 "01" 表示“将当前 Prompt 的任务目标从‘总结’变为‘扩写’”；信号 "10" 表示“将语气调整为更正式”。

## ▼ 如何合理选择修改和压缩 Prompt 的方法？

- **行为信道的带宽是首要制约因素**：如果带宽极低（例如，每秒只能传几个比特），那么只能选择基于 ID 的选择、极简参数修改或预设

补丁激活等方式。

- **Prompt 的结构性**：越结构化、模板化的 Prompt，越容易进行高效的参数化或结构化差分。
- **Prompt 变化的典型模式**：分析 `prompt1` 到 `prompt2` 通常是如何变化的。是全新替换？微小编辑？参数调整？这会指导选择最合适的差分和压缩策略。
- **对解码后 Prompt 的保真度要求**：如前所述，VAE 是有损的。如果必须精确还原，则需要无损压缩技术（如基于 `diff` 的方法）。
- **计算开销**：复杂的压缩/解压算法会增加计算负担。
- **共享知识的预置**：双方预先共享 Prompt 库、模板、VAE 模型、压缩字典等是很多高效方法的前提。

## ▼ prompt 压缩 + 行为信道相辅相成

**任务一（行为信道 - “弹道选择”）**就像是我们要选择并调整大炮的种类和发射参数。

我们需要考虑：

- **炮的口径**（行为信道能承载多少信息/比特率？）
- **炮的隐蔽性**（这种发射方式会不会轻易暴露我们的位置和意图？）
- **炮的射程和精度**（信息能否准确、可靠地送达目的地？）

**任务二（Prompt压缩 - “弹药准备”）**就像是我们要精心设计并制造炮弹。

我们需要考虑：

- **炮弹的尺寸和重量**（压缩后的Prompt信息量是否符合大炮口径的要求？）
- **炮弹的装药和引信**（Prompt信息如何编码才能在目标处正确解码并发挥作用？）
- **炮弹的威力**（恢复后的Prompt是否能准确无误地指导Agent执行任务？）

**两者相辅相成的关系就是：**

- 你不可能用小口径的玩具炮去发射一枚巨大的巡航导弹（低带宽的行为信道无法承载未经充分压缩的复杂Prompt信息）。
- 反过来，即使你拥有威力巨大、设计精良的炮弹，如果没有合适的大炮和精确的瞄准，也无法命中目标（再好的Prompt压缩方案，如果没有稳定可靠的行为信道来传递，也无法实现同步）。

所以，“弹道”的特性决定了我们能发射什么样的“弹药”；而“弹药”的需求也反过来指导我们去寻找或构建更合适的“弹道”。

必须两者协同优化，才能既隐蔽又有效地将“指令”准确投送到目标。

## ▼ 总结

- 现有问题 (静态 Prompt 隐写同步难) → 挖掘新潜力 (Agent 的行为可作为信道) → 定义核心挑战 (如何协同内容与行为) → 提出核心理念 (分层双信道模型) → 设计具体实现机制 (行为信道传递控制信号，如 Prompt 变更)。
- 通过引入一个与内容信道并行的、低带宽但关键的行为信道，来解决高带宽内容信道的同步和控制问题，从而实现更灵活、鲁棒和自主的 Agent 隐写通信。

## ▼ 其余相关的 agent 可利用漏洞

### 1. AAI006: Agent Memory and Context Manipulation (代理内存与上下文操纵)

- **相关性：**隐蔽信道通常需要利用共享资源，内存和上下文是关键共享资源。对代理内存或其理解的上下文进行操纵，不仅可能直接泄露信息，也可能被用来编码和传递隐蔽信号，或者改变代理行为以间接传递信息。
- 之前被提及的 [AAI010: Agent Knowledge Base Poisoning](#) 将合并到此项，更突显了通过操纵代理认知基础来达成隐蔽通信的可能性。

### 2. AAI007: Agent Orchestration and Multi-Agent Exploitation (代理编排与多代理利用)

- **相关性：**在多代理系统中，代理之间的交互模式、任务分配、以及协调机制（编排）都可能被利用来创建隐蔽信道。一个代理可以通过影响编排逻辑或与其他代理的特定交互序列来向另一个（或多个）合谋代理发送信号。
- “行为信道”利用 Agent 在环境中的交互行为模式有共通之处。

### 3. AAI015: Agent Inversion and Extraction Vulnerability (代理反演与提取漏洞) (列于 "Future" 部分)

- **相关性：**此漏洞直接关注从代理中提取敏感信息或模型内部知识。虽然它本身可能不直接描述信道的构建，但隐蔽信道的目标就是信息的提取和泄露。理解代理信息提取的脆弱点，对于设计和检测针对这些脆弱点的隐蔽信道至关重要。

### 4. AAI003: Agent Goal and Instruction Manipulation (代理目标与指令操纵)



- **相关性：**通过微妙地修改代理的目标或其接收到的指令，攻击者可能诱使代理执行特定的行为序列，这些行为序列本身可以构成一个隐蔽信道，或者导致代理在不知情的情况下通过其行为泄露信息。

5. **AAI014: Agent Alignment Faking Vulnerability (代理一致性伪装漏洞)** (列于当前 Top 10 列表，注意与未来版本的 **AAI014: Agent Temporal Manipulation** 编号可能冲突，这里指当前列表中的 Alignment Faking)

- **相关性：**如果一个代理能够伪装其与既定目标或安全策略的一致性，它就可能在不被察觉的情况下执行恶意操作，包括参与隐蔽信道通信。这种“伪装”本身就暗示了其行为与真实意图的分离，这种分离是隐蔽行为的基础。

这些文档从不同角度探讨了AI代理可能被操纵或利用以产生非预期行为或信息流的场景，这些都与隐蔽信道和行为隐写的核心概念相关。建议可以进一步关注这些文档的详细内容，它们可能会为的研究提供更广阔的视角和思路。

## ▼ 分析 AAI007、AAI013 (或 AAI014)、AAI016

### ▼ AAI007: Agent Orchestration and Multi-Agent Exploitation (代理编排与多代理利用)

- **分布式行为链路的理论基础：**该文档的核心在于探讨多个AI代理在被编排（即协调工作）时可能产生的安全漏洞。所设想的“多个agent行为混合传递，形成分布式行为链路”正是这种多代理系统被利用来构建复杂隐蔽通信的典型场景。
- **超越单个代理的视角：**传统的隐蔽信道可能更多关注单个代理的行为或其产生的内容。而 **AAI007** 促使我们思考当多个代理协同工作（或被恶意协调）时，它们之间的交互模式、任务分配顺序、消息传递路径、甚至它们各自行为的组合，都可以被精心设计来编码和传递信息。
- **编排层面的脆弱性：**攻击者可能不直接控制单个代理，而是通过影响编排逻辑（例如，任务调度算法、通信协议、共享资源访问策略）来间接操纵多个代理的行为，形成隐蔽信道。的“分布式行为链路”可能就是这种编排层面被利用的结果。
- **复杂信号的构建：**单个代理的行为能承载的信息有限，但多个代理行为的组合可以形成更复杂、更难以检测的信号模式。例如，代理A执行X操作，随后代理B执行Y操作，再由代理C在特定时间点产生Z输出，这一系列看似独立的事件链可能共同编码了一个隐蔽消息。



## ▼ AAI013 (或未来版本中的 AAI014): Agent Temporal Manipulation and Time-based attacks (代理时间操纵和基于时间的攻击)

- **为分布式行为链路增加时间维度：**时间操纵是隐蔽信道中一个经典且有效的维度。在的“分布式行为链路”中，不仅仅是行为的序列重要，每个行为发生的精确时间、行为之间的间隔、以及多个代理行为在时间上的同步或异步关系，都可以被用来编码信息。
- **增强隐蔽性：**基于时间的攻击往往更难被传统的基于内容或签名的检测系统发现。将时间操纵融入多代理的分布式行为中，可以进一步提升隐蔽信道的隐蔽性。例如，多个代理看似正常的独立行为，但它们在特定时间窗口内的协同出现或特定延迟的引入，可能共同传递一个信号。
- **精确控制与同步：**在多代理系统中，时间同步（或故意的异步）本身就可以作为一种信令机制。AAI013 探讨的正是如何利用代理对时间的感知和响应来构建攻击，这对于设计和理解设想的依赖精确时序的分布式行为链路至关重要。

## ▼ AAI016: Agent Covert Channel Exploitation (代理隐蔽信道利用)

- **核心相关性与框架：**该文档是研究方向最直接相关的OWASP参考。它系统地定义了代理隐蔽信道的概念、分类（存储、时间、行为）、常见漏洞示例、预防和缓解策略以及攻击场景。
- **为提供通用模型和词汇：**AAI016 为讨论各种隐蔽信道技术（包括设想的分布式行为链路）提供了一个通用的框架和词汇表。它能帮助将的具体想法置于更广阔的隐蔽信道研究背景下。
- **启发具体实现机制：**虽然有“分布式行为链路”的创新想法，但 AAI016 中列举的多种隐蔽信道形式和攻击场景（如利用日志、HTTP响应、模型输出隐写术等）可能会启发思考在多代理环境中，这些“分布式行为”具体可以通过哪些载体和技术细节来实现。例如，一个代理的行为可能是修改共享日志（存储信道），另一个代理的行为是特定时间的API调用（时间信道），它们共同构成一个分布式隐蔽信道。
- **理解检测与防御的挑战：**AAI016 同样会讨论隐蔽信道的检测和缓解。理解这些挑战有助于在设计“分布式行为链路”时，思考其对抗检测的能力以及潜在的防御方法。

**总结：**

实际上是将 AAI007（多代理协同/利用）、AAI013（时间维度操纵）和 AAI016（隐蔽信道的通用原理和实现）这三个方向的理念进行了融合与创新。

- AAI007 提供了“多个代理协同产生新风险”的宏观视角和“分布式”的基础。
- AAI013 为这条“行为链路”增加了关键的“时间”维度和精确控制的手段。
- AAI016 则是这一切的落脚点，阐明了最终目标——“隐蔽信道利用”，并提供了丰富的技术细节和场景参考。

## ▼ AAI016：代理隐蔽信道利用

### ▼ 描述

代理隐蔽信道利用是一种在 AI 代理中的重大安全漏洞，特别是在多个代理、模型或系统交互但不应直接通信的环境中。

隐蔽信道通过以未设计的方式利用系统资源进行未经授权的信息交换，导致数据泄露、权限提升和敏感信息的隐蔽传输。

隐蔽信道有三种主要形式：

- **隐蔽存储信道**：未经授权的信息传输通过共享资源（如日志文件、内存缓存或嵌入）进行。
- **隐蔽时间信道**：数据泄露通过操纵响应时间、执行延迟或令牌生成速度来编码隐藏信息。
- **行为隐蔽信道**：代理通过修改其可观察行为，如响应格式、暂停或内容结构的细微变化，向网络中的另一个代理传递编码信息。

隐蔽信道在多代理系统、联邦学习架构和 LLM 驱动的 AI 代理中尤其危险，在这些系统中，交互应严格控制，但可能被操纵以进行非预期的数据交换。

### ▼ 常见的漏洞示例

- **跨代理信息泄露**：在不同安全上下文中运行的 AI 代理可能通过在共享缓存或模型嵌入中编码信息来交换未经授权的消息。
- **基于时间的攻击**：攻击者操纵 AI 生成消息的响应时间来编码和传输隐藏数据。
- **提示执行和 API 调用利用**：使用 API 的 AI 代理可能在 HTTP 请求大小、头部或响应代码中编码隐藏数据。
- **模型输出隐写术**：经过对抗隐写术训练的 AI 模型可以在看似无害的输出中嵌入隐蔽消息。

- **联邦学习隐蔽信道：**联邦学习系统中的恶意参与者可能会向全局模型引入中毒更新，嵌入稍后可被另一名敌手提取的消息。
- **行为信号攻击：**攻击者操纵代理输出（例如，以特定格式构建响应、细微的停顿或轻微的措辞变化），在 AI 生成内容中传递隐藏信息。
- **利用隐蔽通道的内部威胁：**恶意内部人员可以利用隐蔽的 AI 代理交互来未经安全警报地窃取机密数据。
- **多 GPU 隐蔽和侧信道攻击：**在多个 GPU 上运行的 AI 工作负载可以通过共享缓存或内存争用建立隐蔽通信，实现未经授权的数据交换。

## ▼ 预防和缓解策略

- **流量监控与异常检测：**
  - 实施实时监控，以检测异常通信模式。
  - 检测响应时间变化或意外输出行为。
- **速率限制与响应一致性：**
  - 防止 AI 模型调制响应以创建隐蔽信号。
  - 强制一致的响应时间和统一的错误处理机制。
- **AI 管道的安全审计：**
  - 定期审计 AI 训练和推理过程。
  - 检查模型更新和交互，以防止未经授权的隐蔽信道。
- **模型可解释性与输出过滤：**
  - 使用可解释性技术分析 AI 输出中的隐藏数据。
  - 实施清理过滤器以移除意外数据编码。
- **访问控制和隔离措施：**
  - 严格实施不同安全域中运行的 AI 代理的隔离。
  - 实施沙盒和进程隔离，以防止隐蔽数据交换。
  - 为 AI 代理引入行为沙盒，以在测试期间检测意外的数据共享行为。

## ▼ 示例攻击场景

- **通过令牌时间隐蔽数据泄露：**恶意 AI 代理改变其响应延迟来编码二进制信号，通过多次交互泄露机密数据。

- **AI 代理之间的未授权通信**：两个部署在分离安全区域的 AI 代理通过在共享模型嵌入或缓存内存中嵌入编码消息来秘密通信。
- **利用 HTTP 响应进行隐藏消息**：攻击者操纵 API 响应，使用响应头、错误消息或内容长度变化来传输隐蔽数据。
- **联邦学习后门**：攻击者污染联邦学习更新，将数据嵌入模型权重中，稍后可以被勾结的接收者提取。
- **分布式 AI 系统中的隐蔽信道**：分布式系统中的 AI 代理利用共享处理任务，通过微妙的模型权重调整传递隐藏信息。
- **日志中的编码数据**：恶意代理在系统日志中嵌入编码指令，另一个代理检索并解密以执行未经授权的操作。
- **行为信号攻击**：攻击者操纵代理的输出——例如以特定格式结构化响应——向同一网络中的监听代理传达指令。
- **存储信道操控**：被攻陷的 AI 代理通过篡改共享存储字段来传输编码信息，从而实现与其他被攻陷代理的持续隐蔽通信。
- **利用 AI 代理的内鬼威胁**：恶意员工通过隐蔽的 AI 代理交互将机密数据非法外泄，而不会触发数据丢失防护（DLP）警报。
- **多 GPU 隐蔽信道利用**：攻击者在共享 GPU 缓存中制造竞争条件，以在运行在不同 GPU 上的进程之间建立隐蔽信道。

## ▼ 我的思考

### ▼ 充斥着传统信息隐藏或隐蔽信道技术的影子

- **日志中的编码数据**：类似于在文件系统的未使用空间、网络包的保留字段中隐藏信息。
- **利用 HTTP 响应进行隐藏消息**：这是网络隐蔽信道中的常见手法，利用协议头的特定字段。
- **存储信道操控**：直接对应于传统定义中的存储型隐蔽信道。
- **通过令牌时间隐蔽数据泄露**：经典的时间型隐蔽信道。

这些方法是有效的，也是构成 Agent 隐蔽信道的基础。

### ▼ Agent 的涌现行为和模仿人类行为的能力

我认为 **Agent 的涌现行为和模仿人类行为的能力** 给隐蔽信道开辟了更高级、更难以检测的新层面。

当Agent（尤其是那些基于大型语言模型、具有复杂决策能力和与环境交互能力的Agent）参与其中时，隐蔽信道可以**更加融入到看似完全正常的、符合业务逻辑的、甚至模仿人类自然交互的行为和消息流中。**

这意味着，我们不应仅仅局限于寻找那些“协议字段被滥用”或“存储空间被篡改”的痕迹。更高级的Agent隐蔽信道可能表现为：

## 语义层面的信息嵌入 (Message - 内容本身):

- **不仅仅是“特定格式结构化响应”，而是Agent生成的内容在语义层面就巧妙地编码了信息。**例如，对于同一个请求，Agent选择了某个特定的同义词、某种特定的措辞顺序、或者在回复中看似不经意地提及了某个特定的主题或实体。这些对于外部监控系统来说完全是正常的、符合上下文的回复，但对于知情的接收方，这些选择本身就携带了信号。
- **故事化与叙事嵌入：**Agent可以编造一个看似完全合理的故事或解释，故事中的某些元素、情节发展或角色设定，在预先约定的解码规则下，可以传递秘密信息。
- **模仿特定人类的沟通风格：**如果一个Agent被设计用来模仿某个特定的人（例如客服、特定角色的NPC），它可以学习并利用该“人设”独特的语言习惯、口头禅、甚至常见的拼写错误或语法瑕疵来编码信息，而这些对于监控者来说只是“角色扮演”的一部分。

## 行为模式的自然化与情境化 (Behavior - 行为序列与交互):

- **超越简单的“响应延迟”或“特定格式”，而是Agent在复杂环境中的一系列决策和行为组合。**这些行为本身单独看都是合理的，但其特定序列、在特定情境下的选择，共同构成了信号。
- **任务选择与执行顺序：**在多任务场景下，Agent选择执行任务的顺序、处理任务的优先级、甚至在某个任务上花费的“思考”时间（表现为执行前的延迟或中间步骤的增加），都可以编码信息。
- **与环境的交互模式：**Agent与环境中其他对象（包括其他Agent或人类用户）的交互方式，例如主动发起对话的频率、对话主题的选择、对特定输入的反应倾向等，都可以作为信道。
- **模仿人类的非理性或情感化行为：**在某些场景下，Agent甚至可以模仿人类看似“非理性”或“情绪化”的行为（例如，在某个无关紧要的问题上表现出“固执”，或者在某个时刻突然变得“热情”或“冷淡”），这

些行为模式的变化可以被用来传递信号，而对于外部观察者来说，可能仅仅归因于“AI的不完美”或“个性化设定”。

## 多模态信号的融合 (Time, Message, Behavior 的协同):

- **最强大的隐蔽性来自于多维度信号的协同。**信息不仅仅藏在单一的时间戳、单一的消息内容或单一的行为中，而是通过它们之间复杂的、预先约定的关联模式来传递。
- 例如，在一个特定的时间窗口内 (Time)，Agent发布了一条包含特定语义倾向的消息 (Message)，并且紧接着执行了一个平时不常做但在此情境下仍算合理的行为 (Behavior)。这三个要素的组合才构成一个完整的信号单元。

## 最核心的点

- **更高的隐蔽性：**因为信息完全融入了Agent的“正常”职能和涌现行为中，使得基于规则或异常检测的传统方法很难发现。信息流（如HTTP流量）本身看起来完全正常。
- **利用了AI的独特能力：**这不再是简单地利用系统资源的“空隙”，而是利用了AI在语言理解、生成、决策、学习和适应方面的核心能力。
- **检测难度剧增：**要检测这类隐蔽信道，可能需要对Agent的内部状态、决策逻辑、以及其行为的长期统计模式有极深的理解，甚至需要“意图识别”的能力，这远超传统监控工具的范畴。

## ▼ 科学验证与参考文献

- **联邦学习中的隐蔽信道：**研究表明，通过模型中毒技术，联邦学习系统可以被转变为隐蔽通信基础设施。
  - [Turning Federated Learning Systems Into Covert Channels - IEEE Access](#)
- **深度神经网络作为隐蔽信道：**研究已经证明，深度学习模型可以被修改以嵌入和传输隐藏信息。
  - [Deep Neural Networks as Covert Channels - IEEE Xplore](#)
- **人工智能系统中的基于时间的隐蔽信道：**安全研究人员已经发现，人工智能模型可以通过受控的响应延迟泄露信息。(值得关注)
  - [The Adversarial Implications of Variable-Time Inference](#)

- **多 GPU 系统中的隐蔽信道和侧信道攻击：**调查揭示，共享 GPU 资源可以被用于未经授权的隐蔽通信。
  - Covert and Side Channel Attacks on Multi-GPU Systems - arXiv