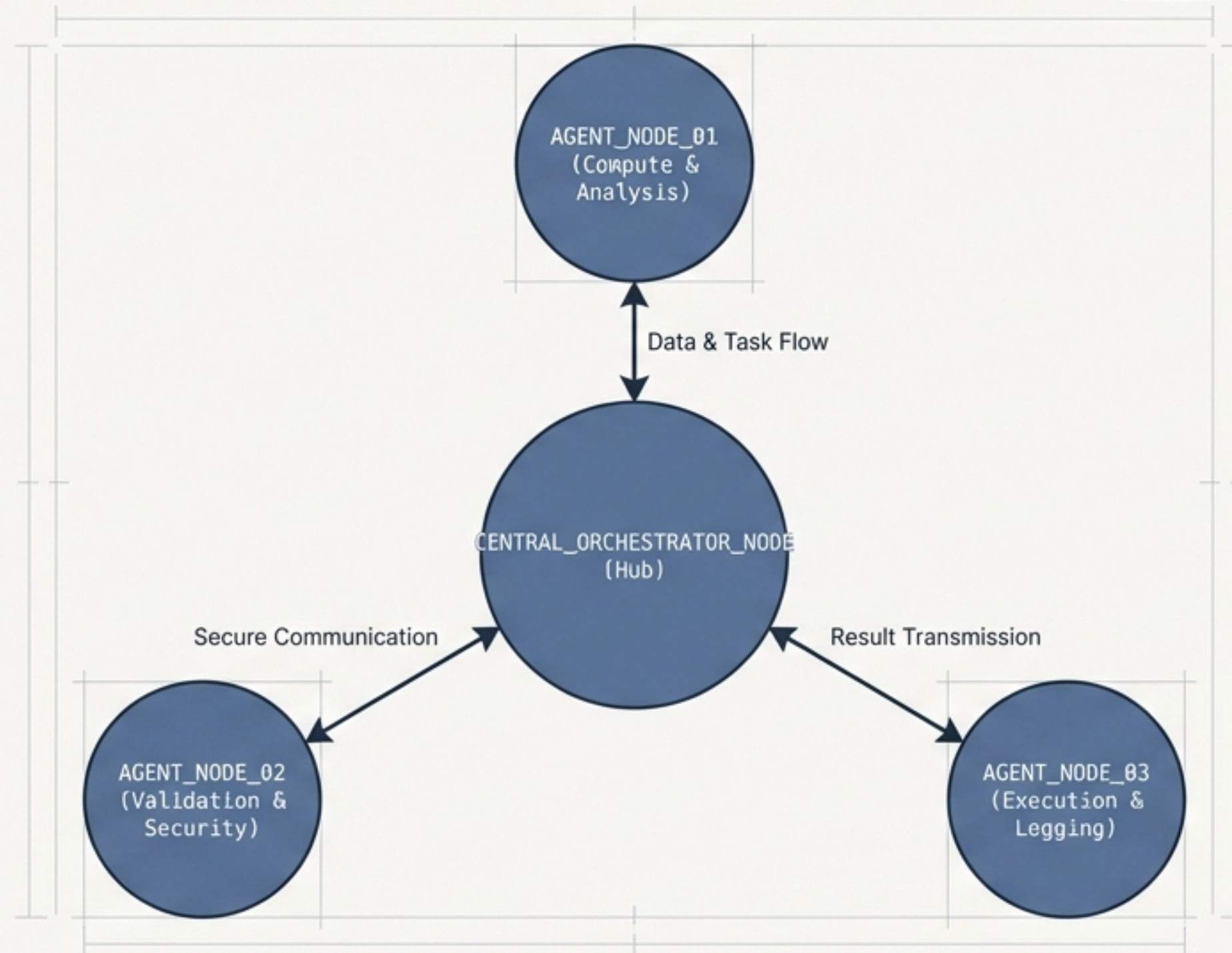


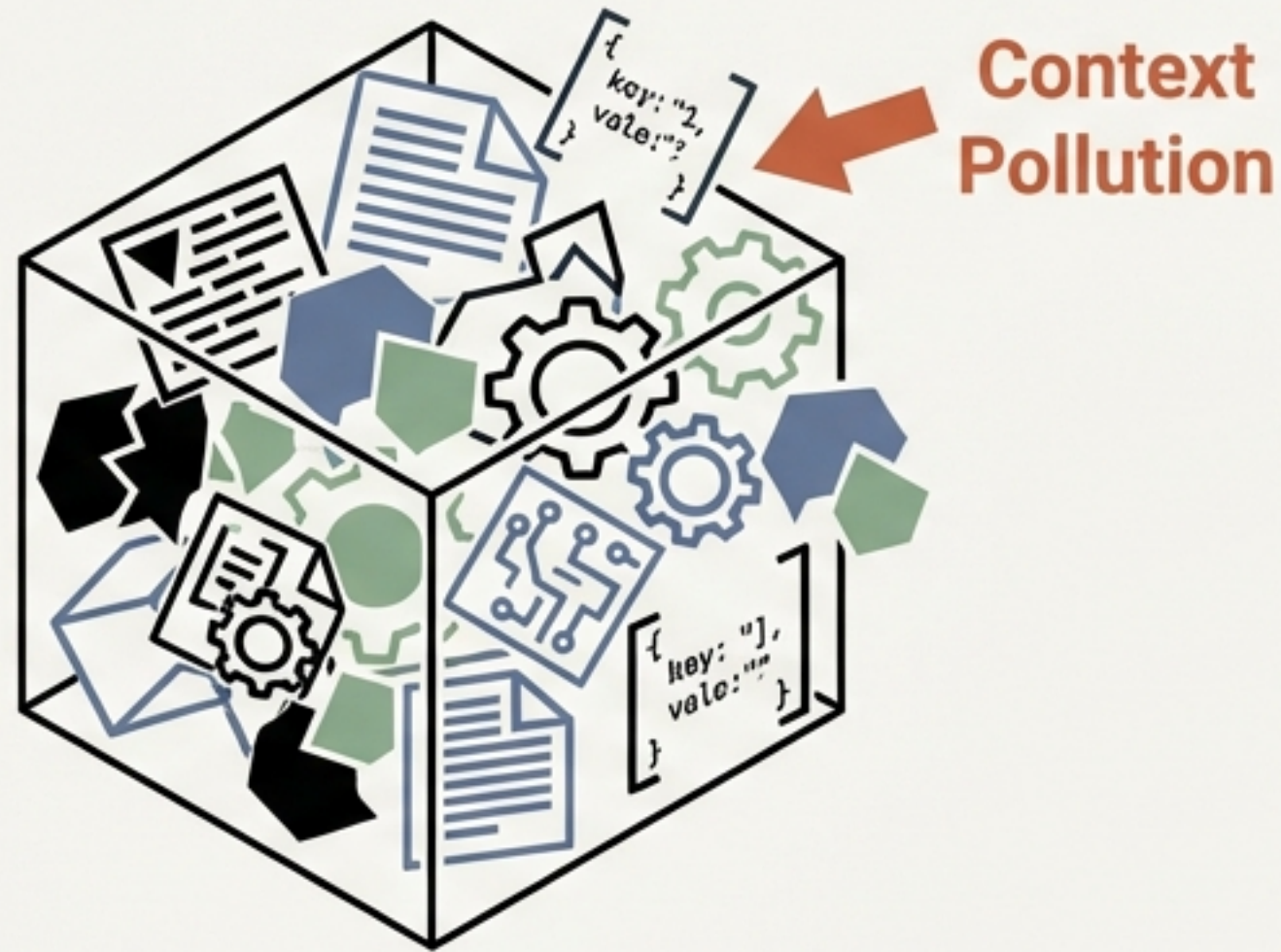
Arquitectura de Orquestación Multi-Agente

Patrones de concurrencia, topologías aisladas y seguridad con `delegate_task`



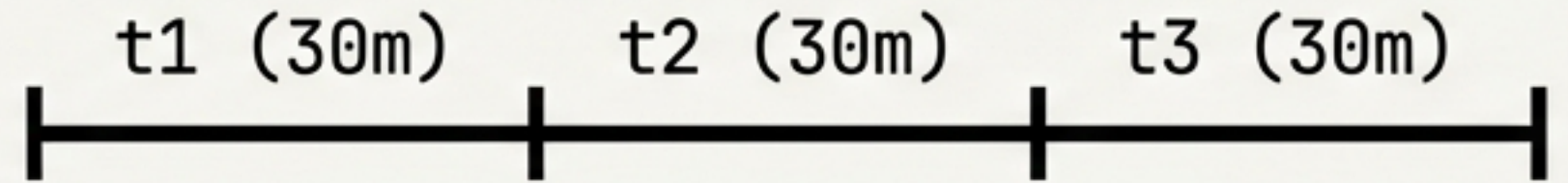
STAR TOPOLOGY: Centralized control point ensures deterministic communication and simplifies failure isolation.

Los Dos Muros del Agente Monolítico



Explosión de Contexto

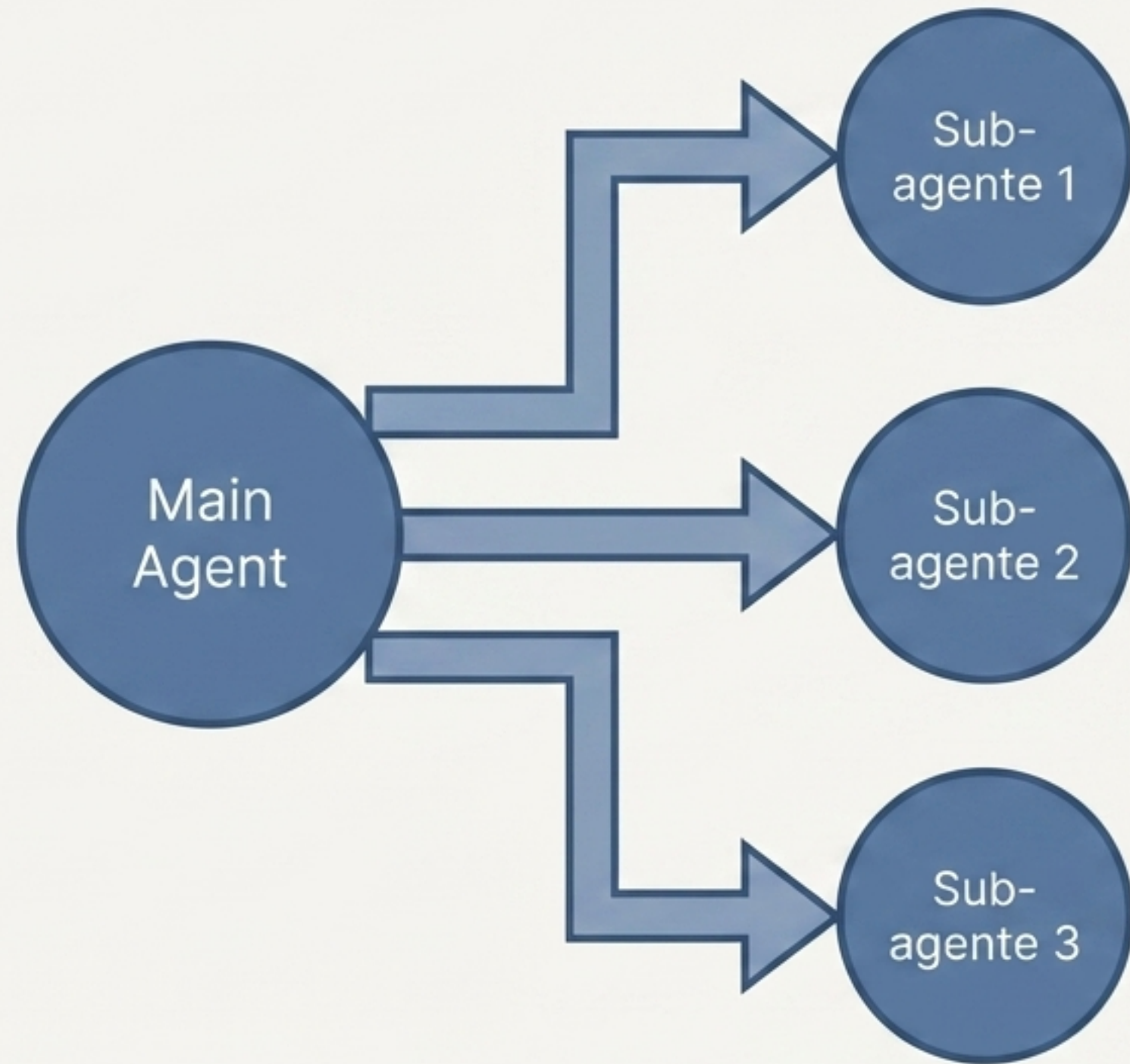
La carga de datos web y de investigación agota la ventana de contexto, dejando espacio insuficiente para el razonamiento lógico de código.



$$\text{Latencia Secuencial} = O(N)$$
$$T_{\text{total}} = t1 + t2 + t3 = 90m$$

Cuello de Botella de Tiempo

La Primitiva Arquitectónica: `delegate_task`



- La herramienta más potente de orquestación nativa.
- Despliega hasta 3 sub-agentes simultáneamente en entornos asilados.

Latencia Paralela = $O(1)$
 $T_{total} = \max(t1, t2, t3)$

Blueprint del Sub-Agente: Aislamiento Estricto

Sub-Agent Sandbox (Memoria Aislada)



**Contexto
Independiente**

Historial de conversación propio.
Evita la contaminación del agente principal.



**Terminal
Aislada**

Sesión independiente a nivel de sistema.
Cero interferencia de procesos.

Retorno de Resultados

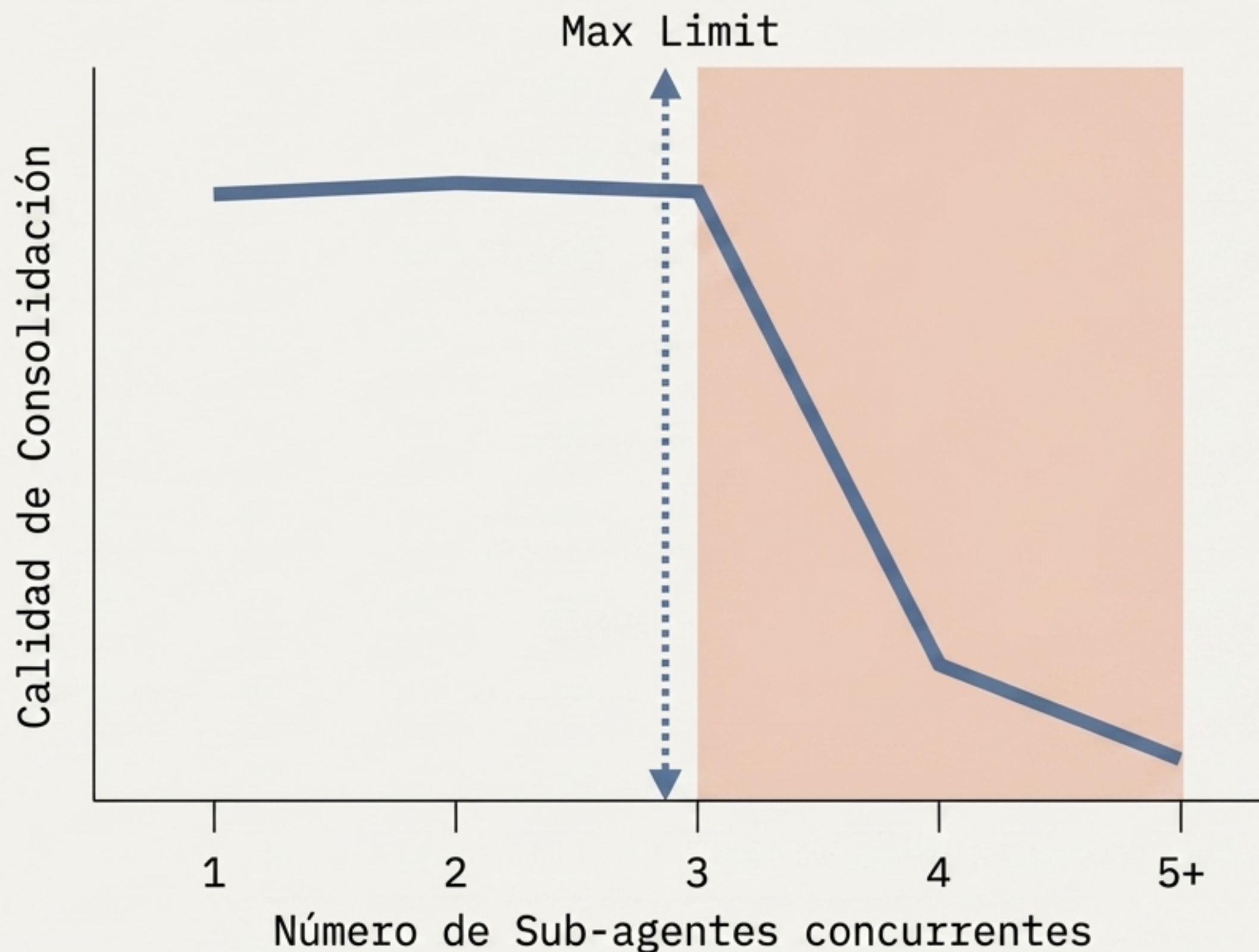


Consolidación

Toolset Restringido (Hard-coded)

[delegate_task, clarify, memory, send_message, execute_code]

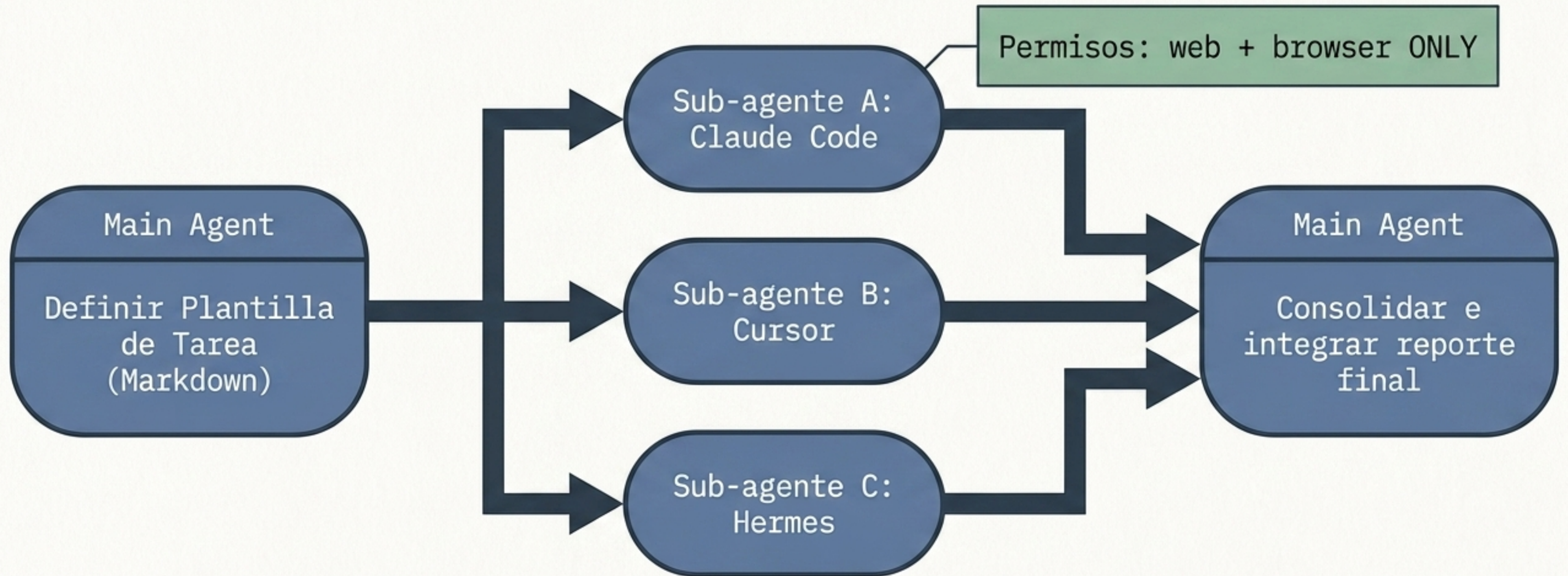
Dispersión de Atención: El Límite Cognitivo



El límite de 3 concurrencias está codificado en el sistema por diseño.

No es una limitación de hardware (compute), es una **prevención estructural** contra la dispersión de atención cuando el LLM consolidador intenta integrar múltiples fuentes independientes.

Flujo de Datos: Análisis Competitivo Paralelo



Impacto en la Latencia de Ejecución

Flujo Secuencial
Monolítico

****Tiempo Total: 40 minutos****

$A + B + C$

Flujo Paralelo
delegate_task

Setup

Task A

Task B

Task C

Consolidate

****Tiempo Total: 15 minutos****

$\max(A, B, C)$

0

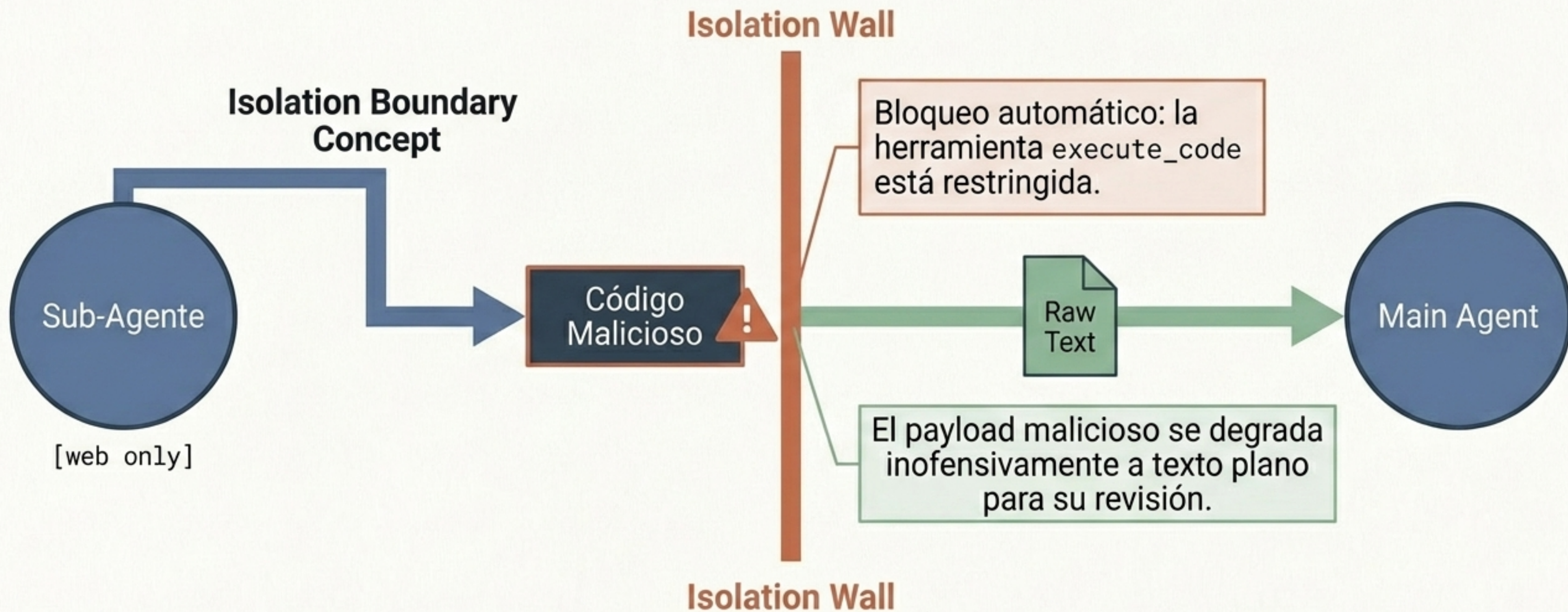
10

20

30

40

Superficie de Ataque y Aislamiento (RCE Mitigation)



El Principio de Menor Privilegio aplicado a nivel de orquestación de agentes.

Matriz de Postura de Permisos

Recomendado (Menor Privilegio)

- **Sub-agentes de Investigación**
→ SOLO [web + browser]
- **Sub-agentes de Código**
→ SOLO [terminal + file + execute_code]
- **Agentes de Consolidación** →
NINGUNA herramienta externa.
Procesamiento de texto nativo.

Antipatrón (Permisos Globales)

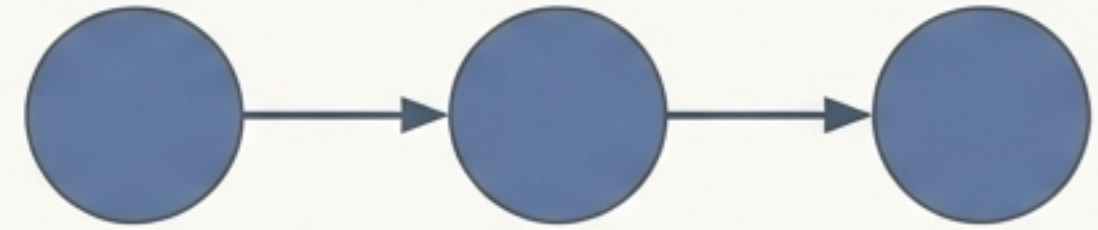
- **Asignación Global** → Cada sub-agente hereda el toolset completo por conveniencia.
- **Impacto Crítico** → Pérdida Pérdida absoluta de la postura de aislamiento seguro. Habilita vectores de ejecución cruzada (Cross-execution risk).

Topologías de Arquitectura: Anthropic vs. Hermes

Anthropic Three-Agent

Dimensiones:

- Roles Fijos (Planificar -> Ejecutar -> Evaluar)
- Secuencial
- Sin memoria nativa

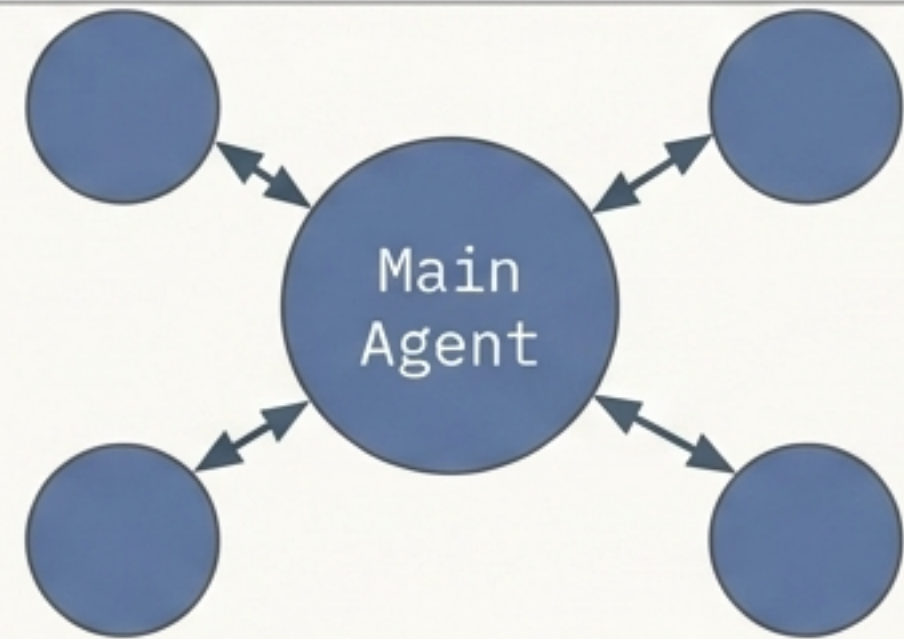


Chain Topology

Hermes delegate_task

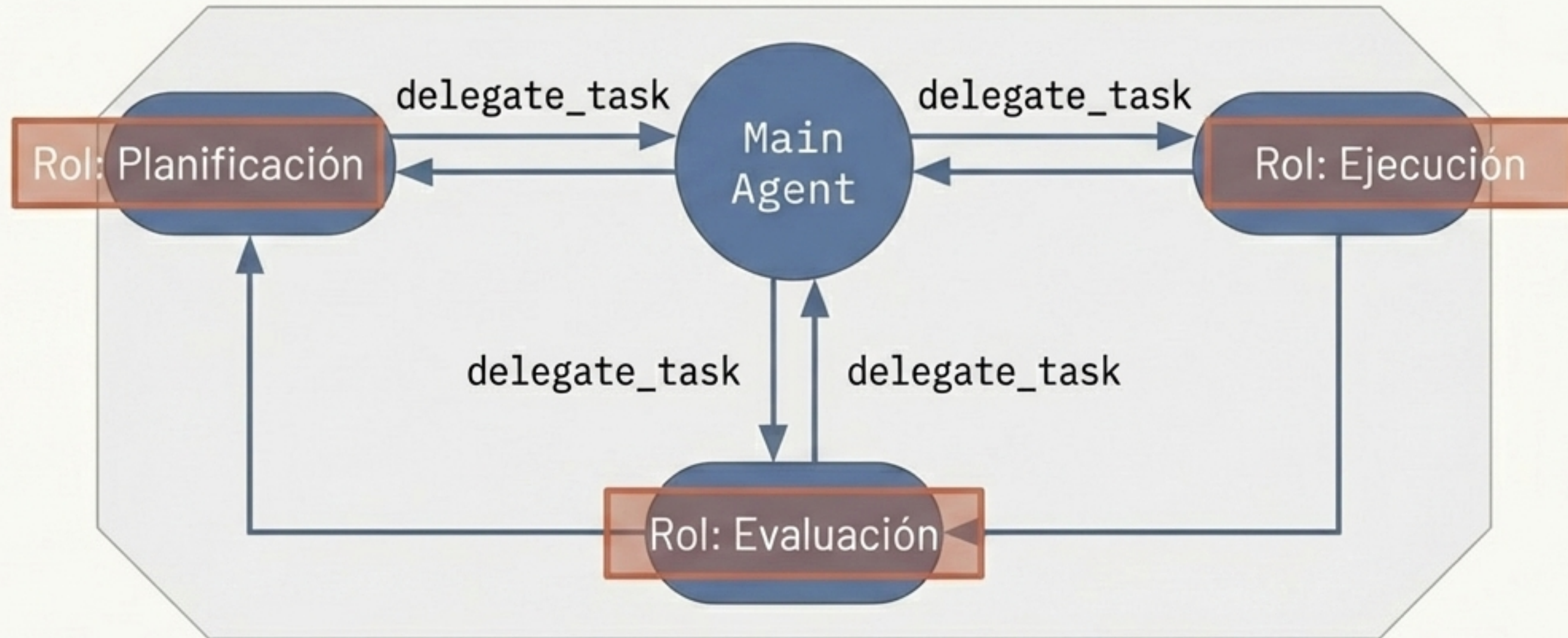
Dimensions:

- Roles Flexibles / Orientados a la tarea
- Hasta 3 hilos concurrentes
- Main Agent retiene memoria global



Star Topology

Diseño Teórico vs. Implementación Física

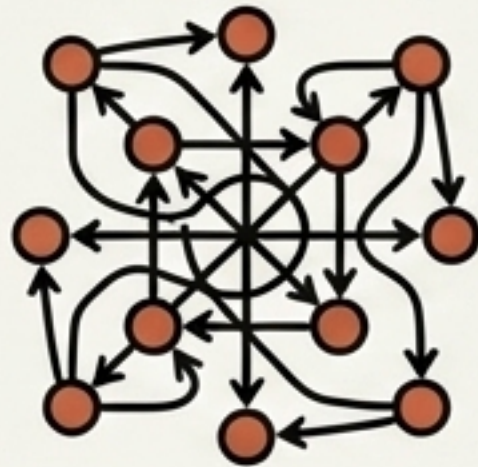


Anthropic gobierna el **DISEÑO**. Es el modelo mental que define cómo descomponer tareas.
`delegate_task` gobierna la **IMPLEMENTACIÓN**. Es la capa de ejecución física que convierte el marco teórico en un ejecutable concurrente.

Heurísticas y Antipatrones de Descomposición

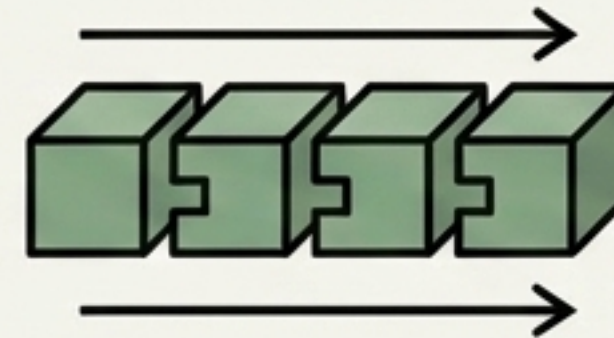
Decomposition Quality Matrix

Antipatrón: Sobre-ingeniería



Si requieres instrucciones de consolidación largas y complejas, la descomposición es incorrecta. Solo usar multi-agente si el contexto es insuficiente o se requiere velocidad paralela pura.

Regla de Oro: Consolidación Simple



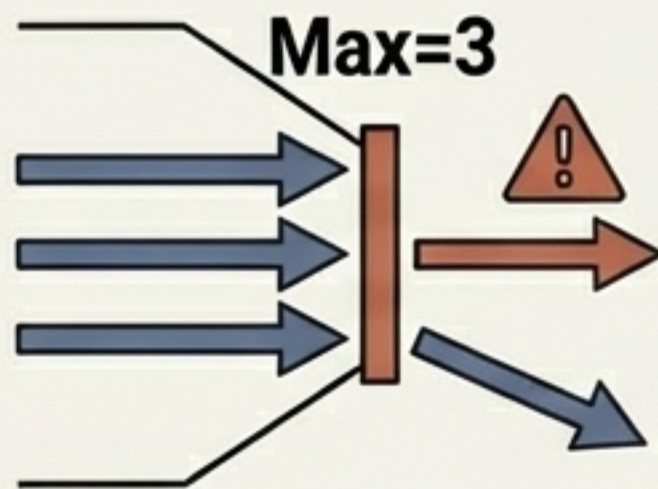
Los outputs del sub-agente deben ser auto-contenidos, con formato uniforme (JSON/Markdown) y directamente componibles sin requerir procesamiento lógico adicional en la capa de retorno.

Blueprint Resumen: Reglas de Orquestación



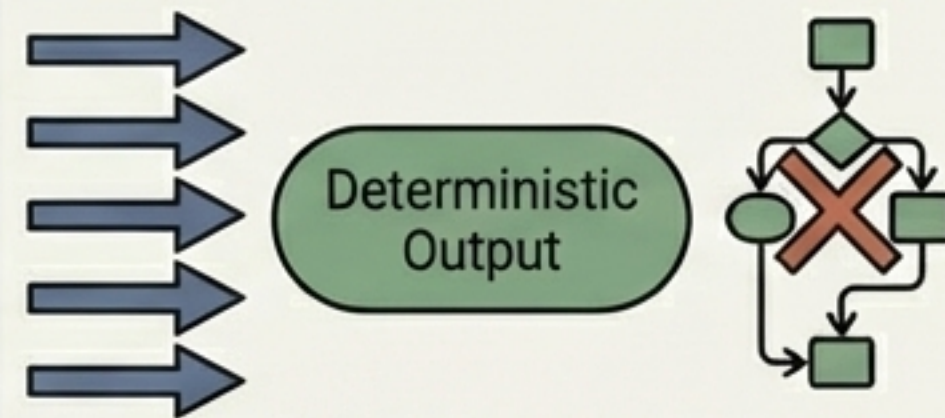
1. Aislamiento por Defecto

Nunca asignar permisos globales. Aislar terminales y aplicar estricto **Principio de Menor Privilegio** mediante restricción de **toolsets**.



2. Topología Límite

Respetar siempre la barrera concurrente **Max=3**. Forzar el límite garantiza un fallo por dispersión de atención en el LLM consolidador.



3. Output Estricto

Diseñar tareas paralelas para retornar resultados **deterministas** y **auto-contenidos**. Evitar evaluación crítica compleja durante la capa de integración.