

The Missing Layer in SOTA AI Controls

A perspective on SOTA AI controls. Anthropic: anthropic.com/engineering/how-we-contain-claude · Google DeepMind: deepmind.google/blog/securing-the-future-of-ai-agents · MACAW: [arXiv:2602.10465](https://arxiv.org/abs/2602.10465) · macawsecurity.com/research

Theme	Anthropic identifies	Google DeepMind identifies	MACAW resolution
Detection ceiling	<u>Probabilistic defenses</u> — non-zero miss rate. Model-layer guardrails shape what Claude “tends to do, not what it can do”	<u>Opaque reasoning ceiling</u> — “As models learn to hide their reasoning... simply reading verbalized chain-of-thought will not be enough”	SecCC extracts <i>effects</i> , not reasoning. <code>credential_read + data_egress</code> → DENY regardless of model internals or obfuscation. The ceiling disappears when you enforce effects, not behavior.
Speed mismatch	<u>User as injection vector</u> — “when the user is typing the instruction, there’s nothing anomalous for a classifier to catch”	<u>Async → sync gap</u> — current system reviews trajectories after the fact; real-time prevention is on the roadmap, not shipped	Every MACAW PEP verification is synchronous — sub-millisecond, inline, before execution. Not a roadmap item. Ships today.
Authorized ≠ safe	<u>Exfiltration through approved domains</u> — destination filter passed an attacker-embedded credential to <code>api.anthropic.com</code>	<u>Insider threat without cryptographic proof</u> — controls are behavioral; a capable insider behaves normally until the moment it acts	Authorization and effect are enforced separately. A valid credential + policy-denied effect = DENY. Every action cryptographically bound to a signed authorization proof.
Alert fatigue at scale	<u>Approval fatigue</u> — users approved ~93% of permission prompts; auto mode catches ~83% of overeager behaviors but ~17% still get through. Any probabilistic defense has a non-zero miss rate	<u>Alert fatigue at scale</u> — most flagged events stem from agent misinterpretation or overeagerness; human review of every trajectory does not scale	Effect-based control is intent-agnostic — whether the agent is misaligned, injected, or overeager, a denied effect combination is blocked automatically. No human review required for the clear-cut cases; policy exceptions are the only escalation path.
Product-scope containment	<u>Containment is product-specific</u> — <code>claude.ai</code> (server container), Claude Code (local sandbox), Cowork (local VM). A poisoned README loads straight into the model’s context despite passing malware checks. No uniform protection across client and server	<u>Internal Google only</u> — roadmap covers GDM internal systems; “model for the wider industry” is aspirational, no external deployment mechanism exists	MACAW protects client and server uniformly — the same cryptographic trust chain secures Claude Code on the developer’s laptop and the Databricks MCP server it calls. Signed invocations, policy enforcement, and tamper-evident audit apply at every boundary crossing regardless of where the agent runs.
Multi-agent trust escalation	<u>Multi-agent trust escalation</u> <i>called out as an open problem — no current solution</i>	<u>Sandboxing without policy composition</u> — no formal policy algebra governs derived agent actions across multi-step workflows	MAPL composition algebra: four proven theorems — monotonic restriction, transitive denial, bounded derivation, no privilege escalation. A sub-agent cannot escalate above its parent’s policy. Not a heuristic.
Persistent memory poisoning	<u>Persistent memory poisoning</u> <i>called out as an open problem — no current solution</i>	<u>No instruction provenance</u> — no mechanism to verify lineage through derived prompt chains; behavioral monitoring is bypassable by a capable model	Authenticated prompts carry cryptographic lineage (parent + root + signatures). Context hash-chains make poisoning detectable by construction, not classification. Provenance is structural.