

RecursiveMAS: Multi-Agent Systems That Think in Latent Space Instead of Text

Kabui, Charles

2026-06-05

[Read at ToKnow.ai](#)

RecursiveMAS: Agents That Think in Latent Space Instead of Text

Latent-state transfer replaces text, making multi-agent AI faster and smarter

+8.3% Average accuracy gain across 9 benchmarks	2.4x End-to-end inference speedup	-75.6% Token usage reduction vs text-based
---	---	--

June 5, 2026

ToKnow.ai

Researchers from UIUC, Stanford, NVIDIA, and MIT have introduced [RecursiveMAS](#), a framework that makes multi-agent AI systems communicate through internal representations instead of natural language. Standard multi-agent setups pass text between agents, meaning each agent must decode its predecessor's full output before generating its own. RecursiveMAS replaces this with a lightweight module called RecursiveLink that transfers latent states directly

between agents, skipping the text bottleneck entirely. The system loops agents in recursive rounds: each pass refines the shared latent state, and only the final agent produces text in the last round. Evaluated across 9 benchmarks spanning math, science, medicine, and code generation, RecursiveMAS averaged +8.3% accuracy over text-based multi-agent baselines while running up to 2.4x faster and using 75.6% fewer tokens. On AIME 2026 math problems, it scored 86.7% compared to 73.3% for the text-based equivalent.

Agents built from small models (1B to 4B parameters each) collectively match or beat single larger models at lower cost. A Planner-Critic-Solver pipeline using three sub-2B models outperformed a single fine-tuned model of similar total size. Because only the RecursiveLink modules are trained (all LLM weights stay frozen), setup requires a fraction of the compute of full fine-tuning. The system supports four collaboration patterns: sequential reasoning, mixture-of-experts, expert-to-learner distillation, and deliberation with tool use.

This work reframes what “multi-agent” means. Previous [research showed](#) that text-based agent debate often fails to beat a single model. RecursiveMAS suggests the problem was never collaboration itself, but the communication medium. Latent-space coordination may be the missing piece that makes multi-agent systems deliver on their promise.

Sources:

- [Recursive Multi-Agent Systems \(arXiv:2604.25917\)](#)
- [RecursiveMAS Project Page](#)
- [RecursiveMAS GitHub Repository](#)
- [RecursiveMAS on HuggingFace Papers](#)

Disclaimer: For information only. Accuracy or completeness not guaranteed. Illegal use prohibited. Not professional advice or solicitation. Read more: [/terms-of-service](#)