

HYPER-MACNET: Task-Aware Hypergraphs for Collective Intelligence in Multi-Agent LLMs

Anonymous ACL submission

Abstract

Multi-agent LLM systems perform well on bounded tasks, but they often break down on cross-domain, long-horizon tasks that require many interdependent constraints to be satisfied together. A key difficulty is coordination: isolated pairwise discussions can yield locally correct plans that become inconsistent when merged. To address this, we present HYPER-MACNET, a multi-agent framework that organizes collaboration with task-aware hypergraphs. Instead of using a fixed communication graph, HYPER-MACNET decomposes each problem into a subtask dependency DAG and assigns each subtask to a hyperedge, which explicitly defines the responsible collaboration unit. It further performs mode-aware collaboration, assigns roles within each unit, and records intermediate artifacts on a global blackboard for dependency-consistent propagation and aggregation. On standard benchmarks, HYPER-MACNET achieves a 6.1% relative gain over the strongest baseline and a 37.8% gain over a vanilla baseline. On the CLM complex-task evaluation, it attains the best mean ranking, indicating more globally consistent coordination under coupled constraints.

1 Introduction

Large language models (LLMs) have demonstrated remarkable proficiency in natural language understanding, code generation, and structured reasoning (Brown et al., 2020; Wei et al., 2022; Achiam et al., 2023). These advances have catalyzed the development of multi-agent systems (LLM-MAS), where specialized agents—assigned heterogeneous roles and tools—collaborate to solve complex problems (Wu et al., 2024; Shinn et al., 2023; Park et al., 2023). However, current evaluations of LLM-MAS predominantly focus on well-scoped tasks with stable reasoning trajectories, such as math and programming benchmarks (Cobbe et al., 2021; Hendrycks et al., 2020; Chen, 2021). In

contrast, real-world deployments demand cross-domain, multi-stage coordination where interdependent constraints must be satisfied across long horizons and intermediate artifacts are continuously revised. In these settings, unstructured natural language interaction often leads to coordination drift, redundant deliberation, and a breakdown in global consistency (Cemri et al., 2025; Valmeekam et al., 2023; Shumailov et al., 2023).

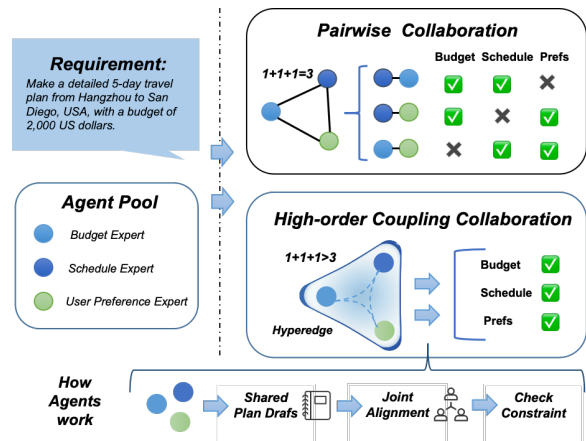


Figure 1: Pairwise exchanges may satisfy constraints only locally and can become inconsistent after post-hoc merging. In contrast, a higher-order collaboration unit co-edits a shared plan draft and performs joint constraint checking over budget, schedule, and preference constraints, producing a feasible travel plan.

We identify a fundamental bottleneck in cross-domain collaboration: *irreducibility* (Figure 1). In many tasks, feasibility depends on jointly aligning heterogeneous constraints under a shared context, and post-hoc merging of independent pairwise threads is therefore often unreliable (Battiston et al., 2021; Abella et al., 2023; Njougou et al., 2025). This aligns with findings in complex systems that higher-order interactions can exhibit emergent effects beyond pairwise links (Battiston et al., 2020; Benson et al., 2016). Consequently, modeling collaboration purely as a standard graph of binary

064	edges can obscure the group-level coupling re-		
065	quired for coherent reasoning. Operationally, ir-		
066	reducibility manifests as a “divide-and-reconcile”		
067	failure mode: constraints may be satisfied in iso-		
068	lated pairwise threads yet become mutually incom-		
069	patible when merged, making aggregation underde-		
070	termined because conflicts surface only under joint		
071	constraint checking.		
072	Hypergraphs provide a principled abstraction for		
073	such higher-order interactions: a hyperedge con-		
074	nects an arbitrary subset of nodes, naturally repre-		
075	senting a discrete collaboration unit (Zhou et al.,		
076	2006; Feng et al., 2019; Yadati et al., 2019). While		
077	recent work uses hypergraphs mainly as communi-		
078	cation backbones to improve message passing effi-		
079	ciency (Zhang et al., 2025), the task structure and		
080	collaboration logic are often left implicit—encoded		
081	in prompts or ad-hoc scripts. Without an explicit		
082	organizational interface, it remains unclear which		
083	subproblems require collective reasoning, how they		
084	depend on each other, and which interaction proto-		
085	col each group should follow (Wang et al., 2024;		
086	Guo et al., 2024).		
087	To bridge this gap, we present HYPER-		
088	MACNET, which uses hypergraphs as an explicit		
089	organizational layer—who collaborates on which		
090	subtask and in what order—rather than only a com-		
091	munication topology. HYPER-MACNET induces		
092	a dependency DAG of interdependent subtasks, and		
093	binds each subtask to a hyperedge (a collaboration		
094	unit) to preserve higher-order coupling throughout		
095	execution. On top of this structure, it selects a col-		
096	laboration protocol and assigns roles within each		
097	team based on subtask semantics, and propagates		
098	intermediate artifacts via a global blackboard along		
099	the DAG to ensure traceability and global consis-		
100	tency. Notably, HYPER-MACNET is parameter-		
101	-free and operates at the system level, enabling plug-		
102	-and-play transfer across LLM backbones without		
103	retraining.		
104	We evaluate HYPER-MACNET on standard		
105	benchmarks (Cobbe et al., 2021; Hendrycks et al.,		
106	2020; Chen, 2021) and a new cross-domain multi-		
107	-step benchmark, CLM. It achieves a 6.1% relative		
108	gain over the strongest baseline and a 37.8% gain		
109	over a vanilla baseline on standard tasks, while		
110	securing the best mean rank on CLM. Our main		
111	contributions are:		
112	• Organizational hypergraph representation.		
113	We formalize hypergraphs to model higher-		
114	-order collaboration under irreducible cou-		
	pling, treating collaboration units (hyper-		115
	edges) as first-class objects.		116
	• Structure induction and protocol routing.		117
	We introduce mechanisms to make subtask		118
	dependencies and interaction procedures ex-		119
	PLICIT, enabling task-adaptive coordination		120
	policies.		121
	• Empirical evaluation. We provide compre-		122
	hensive results on CLM and standard bench-		123
	marks, along with system-level analyses char-		124
	acterizing how structural explicitness affects		125
	coordination quality.		126
	2 Related Work		127
	We briefly review related research from individual		128
	reasoning paradigms to multi-agent collaboration,		129
	high-order topology, and collective intelligence.		130
	Individual Reasoning Paradigms. LLM reason-		131
	ing has evolved from linear CoT to tree/graph struc-		132
	tures such as ToT and GoT (Wei et al., 2022; Yao		133
	et al., 2023; Besta et al., 2024). However, these		134
	paradigms often suffer from structural flattening in		135
	multi-agent settings (Park et al., 2023; Li et al.,		136
	2023): complex internal reasoning is collapsed		137
	into plain-text dialogues, losing topological pri-		138
	ors. Autonomous agents such as AutoGPT and		139
	AutoGen (Richards, 2023; Wu et al., 2024) imple-		140
	ment iterative self-refinement, but still emphasize		141
	individual loops rather than structured collective		142
	synergy.		143
	Coordination and Scaling. Current LLM-based		144
	MAS largely focus on SOP-style workflows or role-		145
	-based organizations or on dynamic group schedul-		146
	-ing and collaboration infrastructures (Hong et al.,		147
	2023; Qian et al., 2024a; Li et al., 2023; Chen		148
	et al., 2024). Recent work on scaling collabora-		149
	-tion with graph-structured topologies (Zhuge et al.,		150
	2024; Qian et al., 2024b; Zhou et al., 2025; Ye		151
	et al., 2025) reveals performance plateaus in dense		152
	dyadic networks. Graph-optimized swarms such		153
	as GPTSwarm and MacNet (Zhuge et al., 2024;		154
	Qian et al., 2024b) mitigate some routing issues		155
	but remain restricted to pairwise links. HYPER-		156
	MACNET instead uses hypergraphs to model irre-		157
	-ducible high-order couplings that go beyond stan-		158
	-dard graphs.		159
	Hypergraphs and Collective Intelligence. Hy-		160
	pergraphs offer a mathematical foundation for		161
	group interactions beyond pairwise links (Battiston		162

et al., 2020). Hypergraph-based communication and forecasting frameworks demonstrate advantages for multi-agent coordination (Zhu et al., 2024; Tian et al., 2022). Existing LLM- or agent-based hypergraph approaches (Zhang et al., 2025) typically treat hypergraphs as static communication backbones. In contrast, we treat hypergraphs as an explicit organizational representation coupled with protocol routing, moving from unconstrained semantic interactions such as debate or voting (Du et al., 2023; Chan et al., 2023) toward a structurally grounded paradigm aimed at globally coherent, super-additive collective intelligence (Woolley et al., 2010).

3 Hypergraph Construction and Structural Alignment

To make higher-order collaboration priors explicit and executable in an LLM-based multi-agent system, we model the organizational structure induced by a task as a weighted hypergraph together with a DAG-dependency hyperedge. In the main text, we focus on the minimal structural definitions required by our protocol routing and cross-hyperedge aggregation modules; a detailed treatment is deferred to the appendix B.

3.1 Formal Definition

We define a weighted hypergraph

$$\mathcal{H} = (\mathcal{V}, \mathcal{E}, w),$$

where $\mathcal{V} = \{v_i\}_{i=1}^N$ is the node set, $\mathcal{E} = \{e_k\}_{k=1}^M$ is the hyperedge set, and each hyperedge $e_k \subseteq \mathcal{V}$ satisfies $|e_k| \geq 2$. The weight function $w : \mathcal{E} \rightarrow \mathbb{R}_{>0}$ assigns a positive weight $w_k = w(e_k)$ to each hyperedge, capturing structural priors such as importance, uncertainty, or difficulty. The node-hyperedge incidence relation can be encoded by a binary incidence matrix indicating whether a node participates in a hyperedge.

3.2 Node-Agent Alignment

Let $\mathcal{A} = \{a_i\}_{i=1}^N$ denote the set of LLM-based agents. We establish a one-to-one correspondence

$$v_i \longleftrightarrow a_i, \quad i \in \{1, \dots, N\},$$

so that \mathcal{V} serves as the structural backbone of the system. All agents share the same base language model, but differ by role prompts and private memory states.

3.3 Hyperedge-Subtask Mapping and Dependency DAG

Let $\mathcal{T}_{\text{sub}} = \{\tau_k\}_{k=1}^M$ be a set of structured subtasks. We define a mapping

$$\psi : \mathcal{E} \rightarrow \mathcal{T}_{\text{sub}}, \quad \psi(e_k) = \tau_k,$$

so that each hyperedge simultaneously specifies who collaborates (its incident agents) and what to solve (the associated subtask and interface constraints). To represent precedence constraints among subtasks, we construct a directed acyclic graph (DAG) over hyperedges,

$$\mathcal{G}_E = (\mathcal{E}, \mathcal{F}),$$

where $(e_k, e_\ell) \in \mathcal{F}$ indicates that τ_k must be completed before τ_ℓ . Since \mathcal{G}_E is acyclic, it admits a topological ordering σ , which constrains macro-level scheduling and aggregation. Dependencies can be compiled from logical/information/verification constraints; in structured settings (e.g., region-, module-, or phase-based decompositions), external structural relations may further induce edges in \mathcal{G}_E .

Overall, Module 1 in Section 4 realizes the mapping

$$(q, \mathcal{A}) \mapsto (\mathcal{H}, \mathcal{G}_E, \psi),$$

providing an executable organizational interface for downstream hyperedge-level collaboration and globally coherent aggregation.

4 Hyper-MACNET Architecture

Building on the task-driven hypergraph formulation in Section 3.1 and the dependency DAG in Section 3.3, we introduce HYPER-MACNET, a higher-order collaboration architecture for LLM-based multi-agent systems. HYPER-MACNET treats each hyperedge as an irreducible collaboration unit (a *team*) and uses the hyperedge dependency DAG to schedule execution and aggregate intermediate results into a globally coherent solution.

HYPER-MACNET consists of three coupled modules: (i) task decomposition and hypergraph induction; (ii) expert collaboration on each hyperedge; (iii) cross-hyperedge aggregation and global decision-making. We additionally discuss an optional extension for structural/weight adaptation, which is disabled in the main experiments.

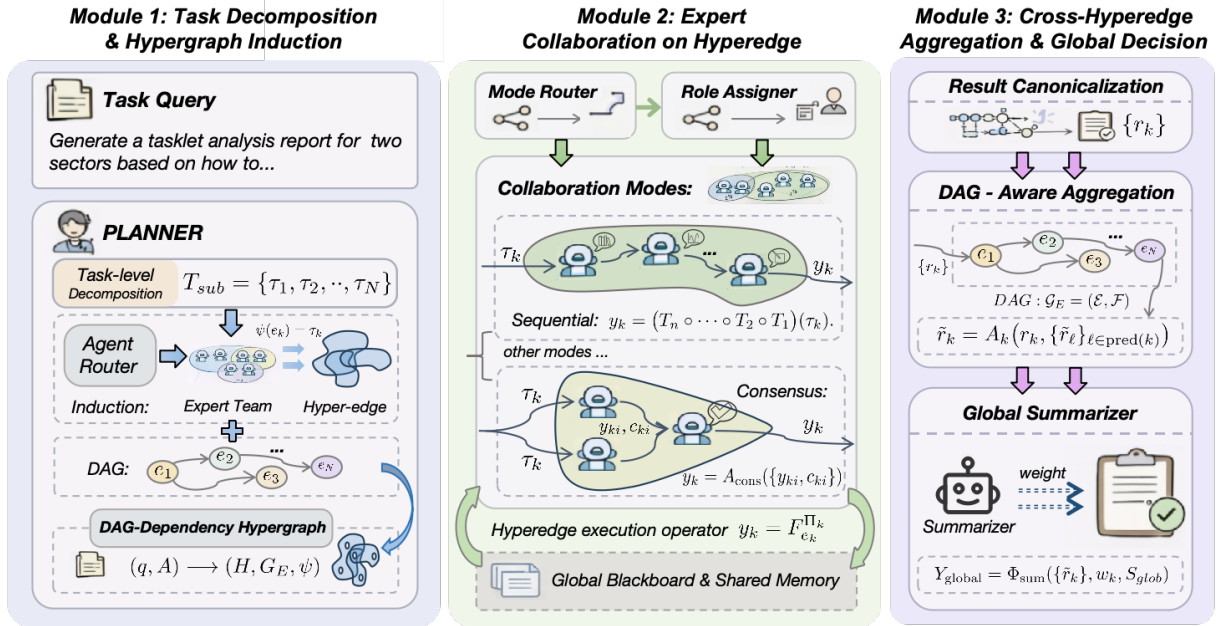


Figure 2: **HYPHER-MACNET**, a hypergraph-structured multi-agent collaboration framework. Module 1 decomposes the query into interdependent subtasks and induces a DAG-dependency hypergraph. Module 2 binds each subtask to a hyperedge and executes it with mode routing and intra-unit role assignment. Module 3 canonicalizes and aggregates hyperedge results along the hyperedge-level dependency DAG to yield a traceable and globally consistent solution. Notation is simplified for illustration; see Sec. 4 for formal definitions.

252 4.1 Module 1: Task Decomposition and 253 Hypergraph Induction

254 Given an input problem q (optionally with external
255 knowledge \mathcal{K}), a planner agent decomposes q into
256 a set of structured subtasks $\mathcal{T}_{sub} = \{\tau_k\}_{k=1}^M$. Each
257 τ_k specifies (i) a concise subproblem description,
258 (ii) required skills/tags for expert routing, and (iii)
259 an expected output interface (e.g., code, structured
260 arguments, statistics, or formatted text).

261 An expert router then selects a collaborating
262 team $\mathcal{A}_k \subseteq \mathcal{A}$ for each τ_k . Using the node-agent
263 alignment in Section 3.2, each team induces a hyper-
264 edge $e_k \subseteq \mathcal{V}$, yielding the hyperedge-subtask
265 mapping $\psi(e_k) = \tau_k$. A positive hyperedge weight
266 w_k is assigned to encode structural priors such as
267 importance, uncertainty, or expected difficulty.

268 Finally, the planner constructs a DAG-
269 Dependency Hypergraph by coupling the induced
270 hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E}, w)$ with a directed acyclic
271 dependency graph over hyperedges, $\mathcal{G}_E = (\mathcal{E}, \mathcal{F})$,
272 where $(e_\ell, e_k) \in \mathcal{F}$ specifies that the subtask
273 $\tau_\ell = \psi(e_\ell)$ must be completed before $\tau_k = \psi(e_k)$.
274 We enforce \mathcal{G}_E to be acyclic, so that it admits
275 a topological order for schedule and aggregate
276 subtasks in precedence order in Module 3.

277 4.2 Module 2: Expert Collaboration on 278 Hyperedges

279 Given the task-driven hypergraph \mathcal{H} and depen-
280 dency graph \mathcal{G}_E , Module 2 specifies how agents
281 collaborate within each hyperedge to solve its asso-
282 ciated subtask.

283 For each hyperedge e_k with team \mathcal{A}_k and subtask
284 τ_k , we instantiate a collaboration protocol Π_k and
285 execute a hyperedge-level operator:

$$286 y_k, \{s_i^+\}_{a_i \in \mathcal{A}_k} = \mathcal{F}_{e_k}^{\Pi_k} \left(\{s_i\}_{a_i \in \mathcal{A}_k}, \tau_k, S_{glob} \right),$$

287 where s_i and s_i^+ denote the pre-/post-collaboration
288 local states of agent a_i , S_{glob} is a shared global
289 state, and y_k is a structured intermediate result.

290 **Mode and role routing.** Mode routing is formu-
291 lated as discrete policy selection over a small, pre-
292 defined set of collaboration modes \mathcal{M} with distinct
293 coordination semantics. We keep the router deter-
294 ministic and parameter-free for reproducibility:

$$295 m_k = \pi_{mode}(\tau_k, e_k, w_k) \in \mathcal{M}.$$

296 Concretely, π_{mode} selects a mode based on
297 lightweight subtask signals (e.g., whether con-
298 straints are tightly coupled or the subtask admits
299 staged refinement): we route to **Consensus** (or

Leader/Decider) when τ_k requires joint reconciliation of coupled constraints, and to **Sequential** when τ_k admits staged refinement; **Brainstorming** is used when broad candidate exploration is beneficial. Roles are assigned by $\pi_{\text{role}}^{(k)} : \mathcal{A}_k \rightarrow \mathcal{R}$, and the protocol Π_k is determined by $(m_k, \pi_{\text{role}}^{(k)})$. We use the following canonical modes:

- **Consensus.** Agents produce candidate drafts and exchange critiques focusing on constraint satisfaction and internal consistency. An integrator merges the candidates:

$$y_k = \mathcal{A}_{\text{cons}} \left(\{y_{k,i}^{(0)}\}_{a_i \in \mathcal{A}_k}, \{c_{k,i}\}_{a_i \in \mathcal{A}_k} \right).$$

- **Sequential.** Agents are ordered into a local pipeline, and each stage applies a transformation to the current artifact. The final output is the composed transformation:

$$y_k = (\mathcal{T}_r \circ \dots \circ \mathcal{T}_2 \circ \mathcal{T}_1)(\tau_k).$$

- **Brainstorming.** A proposer subset generates diverse candidates in parallel, and integrators distill them into a constraint-aligned solution:

$$y_k = \mathcal{A}_{\text{distill}} \left(\{y_{k,i}^{(0)}\}_{i \in \mathcal{P}}, \tau_k \right), \quad \mathcal{P} \subseteq \mathcal{A}_k.$$

- **Leader/Decider.** Team members contribute evidence or partial results, while a designated decider resolves conflicts and commits to a single output:

$$y_k = \mathcal{D} \left(\{y_{k,i}^{(0)}\}_{a_i \in \mathcal{A}_k}, a_{\text{dec}}^{(k)}, \tau_k \right).$$

- **Supervisor/Critic.** Supervisor agents audit intermediate artifacts against task-specific criteria and trigger revisions when violations are detected:

$$y_k = \mathcal{A}_{\text{audit}}(y_{k,\text{draft}}, \text{Rules}(\tau_k)).$$

Module 2 outputs the set of hyperedge-level intermediate results $\mathcal{Y} = \{y_k\}_{k=1}^M$, together with updated local/global states.

4.3 Module 3: Cross-Hyperedge Aggregation and Global Decision

Module 3 integrates hyperedge-level results into a final solution under the dependency constraints encoded by \mathcal{G}_E . The key requirement is dependency-consistent aggregation: each subtask aggregates information from its predecessors.

Canonical result packaging. Each hyperedge output y_k is normalized into a canonical record (a structured text summary with fixed fields, e.g., claims, satisfied/violated constraints, evidence, open issues):

$$r_k = \Phi_{\text{pack}}(y_k, \tau_k, S_{\text{glob}}). \quad (1)$$

Predecessor-consistent aggregation on the DAG. For hyperedge e_k , define its predecessor set

$$\text{pred}(k) = \{\ell \mid (e_\ell, e_k) \in \mathcal{F}\}. \quad (2)$$

We compute dependency-consistent records \tilde{r}_k in a topological order of \mathcal{G}_E :

$$\tilde{r}_k = \mathcal{A}_k \left(r_k, \{\tilde{r}_\ell \mid \ell \in \text{pred}(k)\}, w_k \right). \quad (3)$$

This recursion aligns information flow with the precedence constraints in \mathcal{G}_E , preventing future-information leakage.

Global decision. A global summarizer produces the final answer by explicitly conditioning on the DAG structure, hyperedge weights, and global state:

$$y_{\text{global}} = \Phi_{\text{sum}} \left(\{\tilde{r}_k\}, \mathcal{G}_E, \{w_k\}, S_{\text{glob}} \right). \quad (4)$$

5 Experiments

5.1 Tasks and Datasets

We evaluate HYPER-MACNET on both standard reasoning/code benchmarks. Specifically, we report results on MMLU_hard (a curated subset of more quantitative/formal subjects from MMLU (Hendrycks et al., 2020)), subject list and full model performance on MMLU are shown in Appendix C), GSM8K (Cobbe et al., 2021), and HumanEval (Chen, 2021).

To evaluate long-horizon coordination under coupled constraints, we introduce **CLM** (Complex Long-horizon Multi-task), a dataset of 50 multi-step, cross-domain tasks spanning six MacroSystems. Each task specifies explicit goals and typically includes 3–4 hard, interacting constraints (e.g., physical/temporal, safety/robustness, economic/resource, and ethics/compliance), requiring a structured multi-stage solution or a full reasoning trace; dataset statistics are provided in Appendix E.

CLM is intentionally constructed to isolate irreducible higher-order coupling—the target regime of our method. The claim is not universal dominance, but that when feasibility hinges on joint

Table 1: Comparison of existing single-agent and multi-agent reasoning frameworks and our HYPER-MACNET on standard benchmarks. Columns **Mul.** indicates whether a method uses multiple agents; **Ada.** indicates whether the interaction structure is adaptively organized rather than fixed (✓ = yes, ✗ = no, △ = partial).

Method	Mul.	Ada.	MMLU_hard	GSM8K	HumanEval	Avg.
<i>Single-Agent Methods</i>						
Vanilla	✗	✗	41.65	68.16	43.29	51.03
CoT	✗	✗	<u>52.22</u>	74.26	60.98	62.68
AutoGPT	✗	△	45.68	67.47	48.09	53.75
<i>Static Multi-Agent Topologies (MacNet variants)</i>						
Graph-Chain	✓	✗	44.06	66.26	32.72	47.68
Graph-Star	✓	✗	44.27	67.90	55.49	55.89
Graph-Tree	✓	✗	43.89	69.00	48.78	53.89
Graph-Mesh	✓	✗	43.89	67.40	51.22	54.17
Graph-Layer	✓	✗	43.69	68.08	49.39	53.72
Graph-Random	✓	✗	44.19	67.40	52.44	54.68
<i>Adaptive Multi-Agent Frameworks</i>						
MetaGPT	✓	△	46.54	67.68	64.41	59.28
GPTSwarm	✓	△	47.30	<u>78.43</u>	49.69	58.47
AgentVerse	✓	△	50.06	76.33	72.56	<u>66.32</u>
HYPER-MACNET (Ours)	✓	✓	60.09	79.63	<u>71.34</u>	70.35

constraint alignment, pairwise or unstructured coordination becomes unreliable; CLM operationalizes this setting in a controlled and reproducible way, complementing standard benchmarks.

5.2 Baselines and Implementation Details

Backbone and fairness. Unless otherwise stated, all methods use the same LLM backbone (gpt-3.5-turbo), identical decoding settings (temperature, top-p), and the same total generation budget upper bound per instance: $L=4096$ tokens for standard benchmarks and $L=16384$ tokens for CLM. We disable external tools and retrieval for all methods. For standard benchmarks, we use $N=4$ agents, and for CLM, we use $N=5$ to better match the higher task complexity. All multi-agent methods adopt the same N . We use a unified instruction template and output format across methods to reduce formatting/verbosity biases; baseline prompts and framework configurations are provided in Appendix E.

Baselines. We compare against (i) a single-agent chain-of-thought (Wei et al., 2022) baseline, (ii) AutoGPT-style (Richards, 2023) self-looping single-agent planning, (iii) generic multi-agent frameworks (MetaGPT (Hong et al., 2023), GPTSwarm (Zhuge et al., 2024), AgentVerse (Chen et al., 2024)), and (iv) graph-based multi-agent coordination (MacNet (Qian et al., 2024b)) instantiated with a family of fixed communication topolo-

gies (chain, star, tree, mesh, layer, random). We use Graph-Random as the main MacNet baseline since it provides the best quality-cost compromise among fixed graphs. For each baseline, we only tune a small, shared set of interface-level knobs (role prompt wording, termination rules) for specific tasks under identical trial budgets, without modifying any baseline algorithm.

5.3 Evaluation Metrics and Judging Protocol

Objective benchmarks. For MMLU and MMLU_hard, we report accuracy. For GSM8K, we report answer accuracy under the official evaluation script. For HumanEval, we follow the official evaluation protocol and report pass@1.

CLM: LLM-as-a-judge. CLM outputs are open-ended and multi-criteria; we thus adopt an LLM-as-a-judge protocol. For each instance, we collect each method’s final solution under the same output template, which reduces stylistic and verbosity-driven advantages. For fairness, we anonymize method identities and randomly permute candidates before judging. We use four strong judges (GPT-5.1-think, Gemini-3-pro, Grok-4-reason, DeepSeek-V3.2) under a consistent rubric (cost-effectiveness, feasibility, completeness, final impact, scalability, and alignment) to score and rank all candidates. We report the mean rank aggregated over tasks and judges, with standard de-

Table 2: Ranking of solutions from different models under AI evaluation on the CLM dataset. Each value denotes the mean \pm standard deviation of the ranking score assigned by LLM judges. The same evaluation protocol is used for all compared models.

Method	GPT-5.1-think	Deepseek-V3.2	Gemini-3-pro	Grok-4-reason	Avg.
Vanilla	5.08 \pm 1.68	3.72 \pm 1.62	4.36 \pm 1.35	4.28 \pm 1.51	4.36 \pm 1.54
CoT	5.24 \pm 1.56	6.40 \pm 0.91	6.24 \pm 1.27	5.80 \pm 1.25	5.92 \pm 1.25
MacNet	4.60 \pm 1.85	5.08 \pm 1.47	5.20 \pm 1.73	4.24 \pm 1.98	4.78 \pm 1.76
MetaGPT	4.72 \pm 1.79	4.36 \pm 0.99	3.08 \pm 1.29	3.76 \pm 1.69	3.98 \pm 1.44
GPTSwarm	3.08 \pm 1.19	3.20 \pm 1.50	4.00 \pm 1.38	3.40 \pm 1.70	3.42 \pm 1.44
AgentVerse	4.08 \pm 1.71	3.96 \pm 1.86	3.68 \pm 1.86	4.76 \pm 2.10	4.12 \pm 1.88
HYPER-MACNET (Ours)	1.24 \pm 0.52	1.20 \pm 0.65	1.44 \pm 1.00	1.64 \pm 0.99	1.38 \pm 0.79

Table 3: Performance evaluation across six dimensions on the CLM dataset. Each entry denotes mean of the normalized score assigned by LLM judges. For each dimension, scores are normalized with respect to the Vanilla method, so that Vanilla has value 1.0 and other methods are measured relative to it.

Dimension	CoT	AgentVerse	MacNet	MetaGPT	GPTSwarm	HYPER-MACNET
Cost-effectiveness	-7.42%	0.26%	-5.07%	0.68%	5.96%	14.26%
Feasibility	-7.35%	-1.95%	-4.09%	-0.21%	5.23%	11.86%
Completeness	-13.67%	1.96%	-3.00%	7.62%	7.30%	20.82%
Final impact	-10.29%	1.26%	-4.13%	0.24%	7.79%	19.66%
Scalability	-10.65%	5.53%	-4.06%	8.75%	6.56%	18.01%
Alignment with task	-0.57%	-5.75%	-2.25%	0.20%	7.63%	10.36%

443 viation computed across tasks and repeated eval-
 444 uations. To assess reliability, we repeat judging
 445 $K = 10$ times; ranking stability is high (Kendall’s
 446 $W \in [0.84, 0.87]$ across judges) and score con-
 447 sistency is strong (ICC(3, 10) typically > 0.94
 448 across dimensions). Full prompts, rubrics, aggre-
 449 gation, and reliability analyses are provided in Ap-
 450 pendix D.

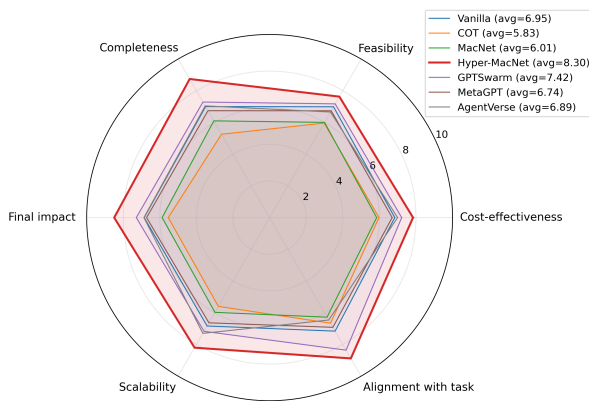


Figure 3: Radar plot of six-dimensional CLM scores, averaged over tasks and LLM judges. HYPER-MACNET achieves consistently highest scores.

5.4 Performance on Complex Tasks (CLM)

451 CLM tasks feature irreducible higher-order cou-
 452 pling: key decisions require jointly aligning het-
 453

454 erogeneous information and interdependent con-
 455 straints within a shared collaboration unit, which
 456 is hard to decompose into independent pairwise ex-
 457 changes and recover by simple aggregation. Under
 458 the LLM-as-a-judge protocol, HYPER-MACNET
 459 achieves the best overall performance, with an aver-
 460 age rank of 1.38 ± 0.79 across tasks and judges, and
 461 it ranks first under all four judges (Table 2). Strong
 462 baselines—generic MAS frameworks (MetaGPT,
 463 GPTSwarm, AgentVerse), and graph-based MAC-
 464 NET (best random fixed-topology variant)—form a
 465 middle group, while the single-agent CoT baseline
 466 performs worst.

467 Judges consistently prefer HYPER-MACNET
 468 because its solutions exhibit clearer multi-stage
 469 structure, stronger constraint coverage, and more
 470 coherent global plans. In contrast, baselines of-
 471 ten spend substantial budget on loosely relevant
 472 branches, yielding lengthy yet partially redundant
 473 deliberation and less stable end-to-end quality.
 474 Even the strongest fixed-graph variant (MACNET)
 475 remains systematically below HYPER-MACNET,
 476 especially on instances requiring explicit multi-
 477 stage planning and cross-domain constraint coordi-
 478 nation. We also observe that pairwise-collaboration
 479 baselines more frequently produce locally consist-
 480 ent partial plans that violate global constraints after
 481 merging, consistent with the proposed divide-and-

Table 4: Ablation study of different components in HYPER-MACNET. We evaluate the contribution of each key component across four benchmarks. The CLM metrics are the average values in Table 3 Δ Avg is the average difference relative to the full HYPER-MACNET.

Methods	CLM	MMLU_hard	GSM8K	HumanEval	Avg.	Δ Avg.
Hyper-MACNET (Full)	15.82	60.09	79.63	71.34	56.72	–
w/o Task DAG (NODAG)	12.26	55.24	77.26	67.07	52.04	-4.68
w/o Hyper-edges (NOHYPEREDGE)	13.30	52.36	75.81	63.41	51.45	-5.28
w/o Mode Switching (SINGLEMODE)	11.68	56.20	77.71	68.90	53.62	-3.10

reconcile failure mode. These results suggest that weakly structured pairwise interactions are insufficient, and that task-DAG-driven hypergraph organization with explicit collaboration modes is crucial for robust multi-agent reasoning on complex tasks.

5.5 Performance on Standard Benchmarks

Table 1 presents a comprehensive comparison on three standard benchmarks. HYPER-MACNET consistently outperforms all baselines, achieving an average score of 70.35 across the three tasks, which constitutes a substantial improvement +4.08 over the strongest baseline AgentVerse (66.32) and representative multi-agent systems such as GPTSwarm and MetaGPT. Concretely, HYPER-MACNET attains 60.09 accuracy on MMLU_hard, 79.63 on GSM8K, and 71.34 on HumanEval. The results highlight the advantage of organizing collaboration through task-aware hypergraphs rather than fixed or pairwise-only communication graphs: HYPER-MACNET more effectively aggregates and coordinates information within collaboration units than static topologies or pipeline-style frameworks. Furthermore, as indicated by the **Mul.** and **Ada.** columns, HYPER-MACNET is among the few methods that simultaneously support genuine multi-agent operation and fully adaptive interaction structures, and this combination translates into the largest performance gains on benchmarks that demand intensive coordination and structured reasoning. Although these benchmarks do not explicitly require long-horizon coordination, task decomposition still induces implicit structure (e.g., reasoning steps or subgoals), allowing hyperedge-level coordination to reduce redundancy and inconsistency even in short-horizon setting.

5.6 Ablation Studies

Table 4 summarizes an ablation study of HYPER-MACNET across four benchmarks. The full model performs best overall (Avg. = 56.72), and removing any component consistently degrades per-

formance, indicating that task-structure induction, higher-order collaboration, and mode-aware coordination are complementary.

Removing the explicit task DAG and topological scheduling decreases Avg. by 4.68 and drops CLM from 15.82 to 12.26, suggesting weaker handling of long-horizon dependencies and cross-step consistency. Replacing hyper-edges with standard pairwise interactions yields the largest degradation (Avg. -5.28), consistent with losing irreducible higher-order coupling: constraints that require joint alignment across multiple agents are more likely to fragment and conflict during aggregation. Disabling mode switching also hurts performance (Avg. -3.10), with a notable CLM drop (15.82 \rightarrow 11.68), reflecting the need to alternate collaboration patterns across stages (decomposition, exploration, consolidation) under a fixed budget.

Overall, the gains stem from jointly modeling task dependencies (DAG), preserving higher-order coupling (hyperedges), and adapting collaboration dynamics via mode routing.

6 Conclusion

In this work, we tackle complex cross-domain reasoning for large language models by introducing HYPER-MACNET, a multi-agent framework that organizes collaboration through task-aware hypergraphs. Instead of relying on a single fixed communication topology, HYPER-MACNET decomposes each problem into a task DAG and uses hyperedges to group role-specialized agents around shared sub-tasks, enabling coherent, structured planning via mode-aware collaboration. Experiments on a cross-domain complex-task and standard benchmarks show that HYPER-MACNET consistently outperforms strong single-agent and multi-agent baselines, highlighting task-structured hypergraph organization as an effective paradigm for scalable multi-agent LLM systems.

562 Limitations

563 **No external tool interfaces.** To isolate the con-
564 tribution of structured collaboration and keep eval-
565 uation controlled and fair, we run all methods with-
566 out external tools (search/retrieval, code execution,
567 solvers, or databases). Future work will study com-
568 position with tool use by enabling plug-in tool call-
569 ing within each hyperedge (e.g., retrieval for fact
570 checking, scripts/solvers for constraint validation)
571 and writing tool outputs back to the shared black-
572 board, alongside executable constraint checkers to
573 improve verifiability and reproducibility.

574 **Text-only setting (no multimodal inputs).** Our
575 study is limited to text-only inputs and does not
576 incorporate visual, auditory, or other multimodal
577 signals. For tasks where constraints, evidence, or
578 feasibility depend on perception (e.g., diagrams,
579 tables, logs, or sensor data), performance remains
580 bounded by pure language reasoning. A natural ex-
581 tension is to encode multimodal evidence as struc-
582 tured constraints or verifiable artifacts on the black-
583 board and integrate them into hyperedge-level co-
584 ordination, enabling agents to jointly reason over
585 language, perception, and tool-derived signals.

586 References

587 David Abella, Piero Birello, Leonardo Di Gaetano,
588 Sara Ghivarello, Narayan G Sabhahit, Christel Siroc-
589 chi, and Juan Fernández-Gracia. 2023. Unravel-
590 ing higher-order dynamics in collaboration networks.
591 *arXiv preprint arXiv:2306.17521*.

592 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
593 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
594 Diogo Almeida, Janko Altenschmidt, Sam Altman,
595 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
596 cal report. *arXiv preprint arXiv:2303.08774*.

597 Federico Battiston, Enrico Amico, Alain Barrat, Gines-
598 tra Bianconi, Guilherme Ferraz de Arruda, Benedetta
599 Franceschiello, Iacopo Iacopini, Sonia Kéfi, Vito La-
600 tora, Yamir Moreno, and 1 others. 2021. The physics
601 of higher-order interactions in complex systems. *Nature physics*, 17(10):1093–1098.

603 Federico Battiston, Giulia Cencetti, Iacopo Iacopini,
604 Vito Latora, Maxime Lucas, Alice Patania, Jean-
605 Gabriel Young, and Giovanni Petri. 2020. Networks
606 beyond pairwise interactions: Structure and dynam-
607 ics. *Physics reports*, 874:1–92.

608 Austin R Benson, David F Gleich, and Jure Leskovec.
609 2016. Higher-order organization of complex net-
610 works. *Science*, 353(6295):163–166.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gersten-
berger, Michal Podstawski, Lukas Gianinazzi, Joanna
Gajda, Tomasz Lehmann, Hubert Niewiadomski, Pi-
otr Nyczyk, and 1 others. 2024. Graph of thoughts:
Solving elaborate problems with large language mod-
els. In *Proceedings of the AAAI conference on artificial
intelligence*, volume 38, pages 17682–17690.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, and 1 others. 2020. Language models are
few-shot learners. *Advances in neural information
processing systems*, 33:1877–1901.

Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A
Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt
Keutzer, Aditya Parameswaran, Dan Klein, Kan-
nan Ramchandran, and 1 others. 2025. Why do
multi-agent llm systems fail? *arXiv preprint
arXiv:2503.13657*.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu,
Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan
Liu. 2023. Chateval: Towards better llm-based eval-
uators through multi-agent debate. *arXiv preprint
arXiv:2308.07201*.

Mark Chen. 2021. Evaluating large language models
trained on code. *arXiv preprint arXiv:2107.03374*.

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang,
Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu,
Yi-Hsin Hung, Chen Qian, and 1 others. 2024. Agent-
verse: Facilitating multi-agent collaboration and ex-
ploring emergent behaviors. In *ICLR*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, and 1 others. 2021. Training verifiers
to solve math word problems. *arXiv preprint
arXiv:2110.14168*.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-
baum, and Igor Mordatch. 2023. Improving factual-
ity and reasoning in language models through multi-
agent debate. In *Forty-first International Conference
on Machine Learning*.

Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong
Ji, and Yue Gao. 2019. Hypergraph neural networks.
In *Proceedings of the AAAI conference on artificial
intelligence*, volume 33, pages 3558–3565.

Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,
Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-
angliang Zhang. 2024. Large language model based
multi-agents: A survey of progress and challenges.
arXiv preprint arXiv:2402.01680.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,
Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
2020. Measuring massive multitask language under-
standing. *arXiv preprint arXiv:2009.03300*.

666	Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, and 1 others. 2023. Metagpt: Meta programming for a multi-agent collaborative framework. In <i>The Twelfth International Conference on Learning Representations</i> .	722
667		723
668		724
669		725
670		726
671		
672		
673	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for "mind" exploration of large language model society. <i>Advances in Neural Information Processing Systems</i> , 36:51991–52008.	727
674		728
675		729
676		730
677		731
678	Thierry Njougouo, Timoteo Carletti, and Elio Tuci. 2025. Collective decision-making with higher-order interactions on d -uniform hypergraphs. <i>arXiv preprint arXiv:2511.13452</i> .	732
679		733
680		734
681		735
682	Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th annual acm symposium on user interface software and technology</i> , pages 1–22.	736
683		737
684		
685		
686		
687		
688	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, and 1 others. 2024a. Chatdev: Communicative agents for software development. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15174–15186.	738
689		739
690		740
691		741
692		742
693		743
694		744
695	Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Kunlun Zhu, Hanchen Xia, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, and 1 others. 2024b. Scaling large language model-based multi-agent collaboration. <i>arXiv preprint arXiv:2406.07155</i> .	745
696		746
697		747
698		748
699		749
700	Toran Bruce Richards. 2023. Autogpt. in https://github.com/significant-gravitas/autogpt , 2023.	750
701		751
702		752
703	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36:8634–8652.	753
704		754
705		755
706		756
707		757
708	Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget. <i>arXiv preprint arXiv:2305.17493</i> .	758
709		759
710		760
711		761
712		762
713	Yu Tian, Xingliang Huang, Ruigang Niu, Hongfeng Yu, Peijin Wang, and Xian Sun. 2022. Hypertron: Explicit social-temporal hypergraph framework for multi-agent forecasting. In <i>IJCAI</i> , pages 1356–1362.	763
714		764
715		765
716		766
717	Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. On the planning abilities of large language models—a critical investigation. <i>Advances in Neural Information Processing Systems</i> , 36:75993–76005.	767
718		768
719		769
720		770
721		771
	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):186345.	772
		773
		774
		775
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	776
		777
	Anita Williams Woolley, Christopher F Chabris, Alex Pentland, Nada Hashmi, and Thomas W Malone. 2010. Evidence for a collective intelligence factor in the performance of human groups. <i>science</i> , 330(6004):686–688.	778
		779
	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, and 1 others. 2024. Autogen: Enabling next-gen llm applications via multi-agent conversations. In <i>First Conference on Language Modeling</i> .	780
		781
		782
	Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, and Partha Talukdar. 2019. Hypergen: A new method for training graph convolutional networks on hypergraphs. <i>Advances in neural information processing systems</i> , 32.	783
		784
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models, 2023. URL https://arxiv.org/abs/2305.10601 , 3:1.	785
		786
	Rui Ye, Shuo Tang, Rui Ge, Yaxin Du, Zhenfei Yin, Siheng Chen, and Jing Shao. 2025. Mas-gpt: Training llms to build llm-based multi-agent systems. <i>arXiv preprint arXiv:2503.03686</i> .	787
		788
	Heng Zhang, Yuling Shi, Xiaodong Gu, Zijian Zhang, Haochen You, Lubin Gan, Yilei Yuan, and Jin Huang. 2025. Hyperagent: Leveraging hypergraphs for topology optimization in multi-agent communication. <i>arXiv preprint arXiv:2510.10611</i> .	789
		790
	Dengyong Zhou, Jiayuan Huang, and Bernhard Schölkopf. 2006. Learning with hypergraphs: Clustering, classification, and embedding. <i>Advances in neural information processing systems</i> , 19.	791
		792
	Heng Zhou, Hejia Geng, Xiangyuan Xue, Li Kang, Yiran Qin, Zhiyong Wang, Zhenfei Yin, and Lei Bai. 2025. Reso: A reward-driven self-organizing llm-based multi-agent system for reasoning tasks. <i>arXiv preprint arXiv:2503.02390</i> .	793
		794
		795
	Tianyu Zhu, Xinli Shi, Xiangping Xu, Jie Gui, and Jinde Cao. 2024. Hypercomm: Hypergraph-based communication in multi-agent reinforcement learning. <i>Neural Networks</i> , 178:106432.	796
		797

776 Mingchen Zhuge, Wenyi Wang, Louis Kirsch,
777 Francesco Faccio, Dmitrii Khizbullin, and Jürgen
778 Schmidhuber. 2024. Gptswarm: Language agents
779 as optimizable graphs. In *Forty-first International
780 Conference on Machine Learning*.

781 A Contents

- 782 • Appendix B: Additional Spectral Hypergraph
783 View and Structural Context
- 784 • Appendix C: Additional Experiments
- 785 • Appendix D: LLM-as-a-Judge Protocol
- 786 • Appendix E: Description and statistical results
787 of the CLM dataset

788 B Additional Spectral Hypergraph View 789 and Structural Context

790 This appendix provides a formal treatment of the
791 incidence structure and normalized operators for
792 weighted hypergraphs. While the main text em-
793 ploys hypergraphs primarily as an executable or-
794 ganizational interface (defining team membership
795 and dependency constraints), the spectral view pre-
796 sented here offers an analyzable structural prior.
797 This formulation serves as an interpretability lens
798 that quantifies the couplings induced by shared
799 hyperedges and supports a lightweight, structure-
800 weighted aggregation mechanism (Section B.6).
801 Crucially, this spectral analysis is intended as an
802 analytical tool rather than a learned component;
803 it informs the system’s coordination logic without
804 introducing any trainable parameters, thereby main-
805 taining the framework’s parameter-free nature.

806 B.1 Incidence, Weights, and Degrees

807 Consider a weighted hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E}, w)$
808 with $|\mathcal{V}| = N$ and $|\mathcal{E}| = M$. The node–hyperedge
809 incidence structure is encoded by a binary inci-
810 dence matrix

$$811 \mathbf{H} \in \{0, 1\}^{N \times M}, \quad \mathbf{H}_{ik} = \begin{cases} 1, & v_i \in e_k, \\ 0, & v_i \notin e_k. \end{cases}$$

812 Let $\mathbf{W} = \text{diag}(w_1, \dots, w_M)$ be the diagonal ma-
813 trix of hyperedge weights with $w_k = w(e_k) > 0$.

814 We define the (weighted) node degree and the
815 hyperedge size (degree) as

$$816 d(v_i) = \sum_{k=1}^M w_k \mathbf{H}_{ik}, \quad \delta(e_k) = \sum_{i=1}^N \mathbf{H}_{ik},$$

817 and the corresponding diagonal degree matrices

$$818 \mathbf{D}_v = \text{diag}(d(v_1), \dots, d(v_N)),$$

$$819 \mathbf{D}_e = \text{diag}(\delta(e_1), \dots, \delta(e_M)).$$

821 Intuitively, \mathbf{W} emphasizes important hyperedges,
822 \mathbf{D}_e^{-1} normalizes for team size, and $\mathbf{D}_v^{-1/2}$ prevents
823 high-participation nodes from dominating propaga-
824 tion.

825 B.2 Normalized Hypergraph Adjacency and 826 Laplacian

827 Following standard constructions in spectral hyper-
828 graph theory, we define the normalized hypergraph
829 adjacency operator

$$830 \mathbf{A}_{\mathcal{H}} = \mathbf{D}_v^{-\frac{1}{2}} \mathbf{H} \mathbf{W} \mathbf{D}_e^{-1} \mathbf{H}^T \mathbf{D}_v^{-\frac{1}{2}} \in \mathbb{R}^{N \times N}, \quad (5)$$

831 and the corresponding normalized hypergraph
832 Laplacian

$$833 \mathbf{L}_{\mathcal{H}} = \mathbf{I} - \mathbf{A}_{\mathcal{H}}. \quad (6)$$

834 The operator in Eq. (5) admits a natural two-stage
835 interpretation: information propagates from nodes
836 to incident hyperedges and back to other nodes
837 within the same hyperedges, with contributions
838 modulated by hyperedge weights and normalized
839 by team sizes and node degrees.

840 B.3 Coupling Strength Induced by Shared 841 Hyperedges

842 We interpret each entry of $\mathbf{A}_{\mathcal{H}}$ as a *structurally*
843 *induced coupling strength*. Define

$$844 \alpha_{ij} = (\mathbf{A}_{\mathcal{H}})_{ij}.$$

845 Expanding Eq. (5) yields the closed-form expres-
846 sion

$$847 \alpha_{ij} = \frac{1}{\sqrt{d(v_i) d(v_j)}} \sum_{k=1}^M \frac{w_k}{\delta(e_k)} \mathbf{H}_{ik} \mathbf{H}_{jk}. \quad (7)$$

848 Thus, α_{ij} increases when v_i and v_j co-occur in
849 many hyperedges, especially high-weight ones, and
850 decreases for large hyperedges (via $\delta(e_k)$) and high-
851 participation nodes (via $d(v_i)$), which improves
852 interpretability and prevents trivial amplification.

853 B.4 Structural Context Operator (Optional)

854 For analysis or optional structural bias, let $x_i^t \in$
855 \mathbb{R}^d denote a vector representation of the state of
856 agent a_i (aligned with node v_i) at step t , e.g., an
857 embedding of a *structured textual state record* (role
858 description, memory summary, and intermediate
859 artifacts). Stacking row-wise gives $\mathbf{X}^t \in \mathbb{R}^{N \times d}$.
860 We define the structural context transformation

$$861 \tilde{\mathbf{X}}^t = \mathbf{A}_{\mathcal{H}} \mathbf{X}^t, \quad \tilde{x}_i^t = \sum_{j=1}^N \alpha_{ij} x_j^t, \quad (8)$$

862 which aggregates neighbor states according to cou-
863 plings in Eq. (7). This operator is not required
864 by our main pipeline and introduces no trainable
865 parameters; we use it mainly as an interpretable
866 structural prior.

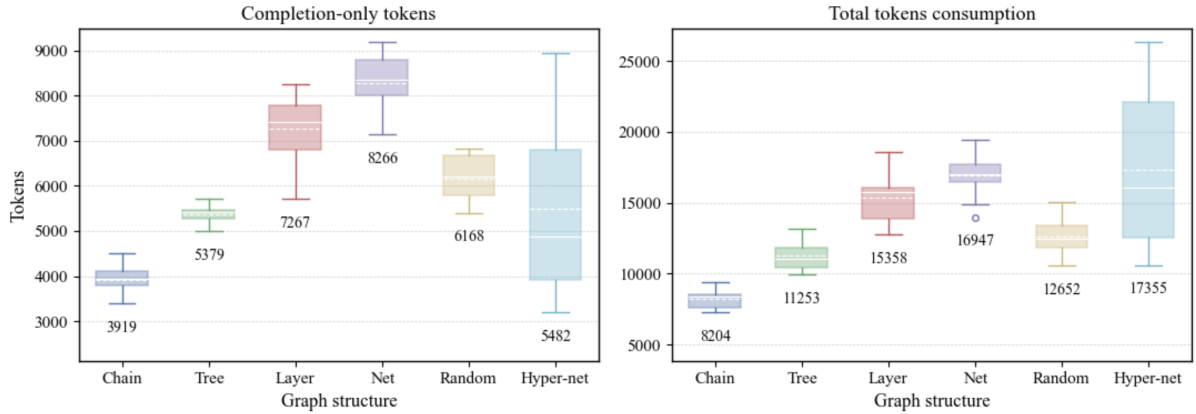


Figure 4: CLM token usage under different collaboration topologies; We vary only the interaction topology, with the left plot showing completion tokens and the right plot showing total tokens. Hyper-net corresponds to our task-aware HYPER-MACNET architecture.

B.5 Irreducibility Beyond Pairwise Graph Expansions

A common reduction from hypergraphs to ordinary graphs is a pairwise (clique) expansion, e.g.,

$$\mathbf{A}_{ij}^{\text{clique}} = \sum_{k=1}^M \mathbb{I}[\mathbf{H}_{ik}\mathbf{H}_{jk} = 1],$$

which replaces each hyperedge by pairwise links. However, such expansions do not faithfully capture (i) team-size normalization (the role of $\delta(e_k)$), (ii) hyperedge importance weighting (the role of w_k), and (iii) the execution semantics of a *joint* collaboration unit (a hyperedge team) as opposed to a collection of independent dyadic interactions. The normalized operator in Eq. (5) preserves these factors through \mathbf{W} , \mathbf{D}_e^{-1} , and $\mathbf{D}_v^{-1/2}$, providing a more principled characterization of higher-order couplings.

B.6 Structure-Weighted Aggregation Interface

To avoid disconnect between the spectral formalization and the system implementation, we describe a minimal interface that can be optionally plugged into cross-hyperedge/global aggregation. Let u_j denote an evidence vector (or scalar confidence) produced by node/agent j (e.g., an embedding of an intermediate artifact, a scored candidate, or a structured record). We define a structure-weighted aggregation for node i as

$$\text{Agg}_i(\{u_j\}_{j=1}^N) = \sum_{j=1}^N \alpha_{ij} u_j, \quad (9)$$

where α_{ij} is given by Eq. (7). This aggregation prioritizes evidence from structurally coupled agents and attenuates signals from weakly related ones, supporting globally coherent integration under higher-order collaboration structure. Importantly, Eq. (9) can be used as a lightweight bias term or an analysis tool and does not require any additional learning beyond the base LLM.

C Additional Experiments

C.1 Inference Cost under Different Collaboration Topologies

To analyze how collaboration topology affects the cost of multi-agent reasoning, we fix the backbone model, the number of agents (five), and all prompt templates, and vary only the interaction structure. On the CLM dataset, we instantiate MacNet with several representative static topologies (Chain, Tree, Layer, Net, Random) and compare them with our proposed HYPER-MACNET. and plot the token distributions in Figure 4.

Overall, the static topologies exhibit a clear trade-off between structural simplicity and communication cost: simpler structures such as Chain and Tree tend to require fewer tokens but offer limited capacity for rich coordination, whereas denser structures such as Layer and Net allow more information exchange but quickly amplify the number of messages and tokens, especially on harder tasks. Random lies in between and provides a reasonable but still task-sensitive quality–cost compromise, which is why we adopt it as the MacNet baseline in our main comparisons.

HYPER-MACNET differs from these static designs in that its “Hyper-net” structure is not a sin-

929	gle fixed graph, but a task-aware hypergraph that	sequences, and analytic geometry. In the hard	977
930	activates different hyper-edges and collaboration	subset, problems often combine multiple con-	978
931	modes according to the task DAG and subtask de-	cepts and require non-trivial logical inference.	979
932	pendencies. As a result, its token usage reflects an		
933	adaptive adjustment of structural complexity: for	• high_school_statistics. Focuses on basic	980
934	relatively straightforward tasks, the induced inter-	probability and statistics, including distribu-	981
935	action pattern remains close to sparse topologies,	tions, expectation/variance, and elementary	982
936	while for constraint-heavy or cross-domain tasks,	hypothesis testing, emphasizing statistical rea-	983
937	local regions of the hypergraph become denser	soning and interpretation.	984
938	where additional coordination is needed. Empiri-		
939	cally, HYPER-MACNET operates in a similar or	• machine_learning. Includes conceptual ques-	985
940	slightly higher cost regime than the stronger static	tions on supervised/unsupervised learning,	986
941	topologies, but achieves significantly better judged	loss functions, optimization, and generaliza-	987
942	quality on CLM and shows more stable token usage	tion, emphasizing theoretical understanding	988
943	across tasks, indicating a better quality-per-token	and method comparison rather than implemen-	989
944	trade-off than hand-tuned fixed graphs.	tation details.	990
945			
946	C.2 MMLU and MMLU_hard	Results of MMLU. On the MMLU benchmark,	991
947	Subjects of MMLU_hard. The MMLU-Hard	Hyper-MACNET achieves the highest overall ac-	992
948	benchmark evaluates model performance on a sub-	curacy at 72.28%, outperforming all baseline meth-	993
949	set of particularly challenging questions drawn	ods by a clear margin (Table 5). Compared with	994
950	from the original MMLU suite. The following	the vanilla single-agent baseline (66.36%), Hyper-	995
951	subjects are included in our evaluation (total 1000	MACNET yields a gain of nearly 6 percentage	996
952	tasks):	points, indicating the effectiveness of structured	997
953		multi-agent coordination beyond simple prompting.	998
954	• abstract_algebra. Covers groups, rings,	While chain-of-thought prompting (CoT), Auto-	999
955	fields, homomorphisms, and quotient struc-	GPT, and other multi-agent frameworks such as	1000
956	tures. Questions emphasize formal definitions,	GPTSwarm and Agentverse provide moderate im-	1001
957	theorem-level reasoning, and structural prop-	provements over the vanilla baseline, their perfor-	1002
958	erties rather than numerical computation.	mance remains clustered around 68–69%, suggest-	1003
959		ing limited benefits from unstructured or fixed inter-	1004
960	• college_computer_science. Includes under-	action patterns. In contrast, Hyper-MACNET con-	1005
961	graduate topics in algorithms and data struc-	sistently surpasses these approaches, highlighting	1006
962	tures, computational complexity, operating	the advantage of explicitly modeling task structure	1007
963	systems, computer architecture, and program-	and higher-order collaboration when addressing	1008
964	ming language principles, focusing on algo-	heterogeneous, multi-domain questions in MMLU.	1009
965	rithmic reasoning and foundational system-		
966	level understanding.	D LLM-as-a-Judge Protocol	1010
967			
968	• college_mathematics. Encompasses standard	Many tasks evaluated in this work are open-ended	1011
969	undergraduate mathematics (e.g., calculus, lin-	and multi-criteria, and thus do not admit a sin-	1012
970	ear algebra, probability theory, differential	gle verifiable ground-truth solution. In such set-	1013
971	equations). Questions typically require multi-	tings, candidate solutions may differ in structure	1014
972	step derivations and careful conceptual rea-	and trade-offs (e.g., feasibility vs. completeness),	1015
973	soning.	making standard automatic metrics insufficient. To	1016
974		enable scalable and reproducible evaluation, we	1017
975	• college_physics. Covers introductory uni-	adopt an LLM-as-a-Judge protocol to score and	1018
976	versity physics such as classical mechanics,	rank solution-level outputs.	1019
	electromagnetism, thermodynamics, and ba-		
	sic modern physics, assessing both qualitative	D.1 Candidate Collection and Judging Models	1020
	physical intuition and quantitative analysis.		
		For each task instance, we collect the final outputs	1021
	• high_school_mathematics. Draws from	produced by all methods. When a task specifies	1022
	secondary-level algebra, geometry, functions,	explicit requirements, we additionally collect the	1023

	Vanilla	CoT	AutoGPT	MacNet-best	GPTSwarm	Agentverse	Hyper-MACNET
Accuracy(%)	66.36	69.00	67.58	68.77	69.11	68.54	72.28

Table 5: Overall accuracy on MMLU.

1024 corresponding structured plans and necessary rea- 1065
1025 soning artifacts, so that judges can verify require- 1066
1026 ment coverage and constraint satisfaction. 1067

1027 We use a panel of high-performance LLMs 1068
1028 as judges: GPT-5.1-think, Gemini-3-pro, 1069
1029 Grok-4-reason, and DeepSeek-V3.2. Each judge 1070
1030 independently scores and ranks all candidates for 1071
1031 each task. 1072

1032 D.2 Anonymization and Randomization

1033 To reduce identity- and presentation-induced bias, 1073
1034 we (i) anonymize all candidate solutions by remov- 1074
1035 ing any source-identifying information, and (ii) ran- 1075
1036 domly shuffle the order of candidates for each task 1076
1037 instance prior to evaluation. All judging models are 1077
1038 prompted with the same unified evaluation prompt 1078
1039 to ensure consistent criteria. 1079

1040 D.3 Scoring Dimensions and Rubric Prompt

1041 Each candidate is evaluated along six dimensions 1080
1042 on a 10-point scale (rounded to two decimals; base- 1081
1043 line score 6.00). We use the following anchor 1082
1044 rubric for calibration. Scores may take intermedi- 1083
1045 ate values (e.g., 7.50) when a solution falls between 1084
1046 anchors. 1085

1047 (1) Cost-effectiveness.

- 1048 • **3-point anchor:** Disproportionately high in- 1086
1049 put cost (compute, interaction, or implementa- 1087
1050 tion) for marginal gains; clear resource waste 1088
1051 or “overkill”. 1089
- 1052 • **6-point anchor:** Costs are broadly propor- 1090
1053 tional to achieved quality; comparable to com- 1091
1054 mon practice; neither notably efficient nor 1092
1055 wasteful.
- 1056 • **9-point anchor:** Achieves strong results with 1093
1057 minimal additional cost; high leverage (large 1094
1058 quality gain per unit cost); avoids unnecessary 1095
1059 overhead. 1096

1060 (2) Feasibility.

- 1061 • **3-point anchor:** Relies on unavailable capa- 1097
1062 bilities, unrealistic assumptions, or extremely 1098
1063 complex execution steps; high risk of failure 1099
1064 in practice. 1100
1101

- **6-point anchor:** Logically sound and imple- 1065
mentable, but requires non-trivial expertise, 1066
engineering effort, or a standard developmen- 1067
t/execution cycle. 1068

- **9-point anchor:** Simple, actionable, and near 1069
plug-and-play; clear steps with low barrier 1070
to execution; minimal dependency on special 1071
resources. 1072

(3) Completeness.

- **3-point anchor:** Missing essential steps; un- 1074
clear or absent input/output specification; ig- 1075
nores boundary conditions and potential fail- 1076
ure modes. 1077

- **6-point anchor:** Covers the main workflow 1078
and solves the nominal case, but has gaps 1079
for edge cases, exceptions, or secondary con- 1080
straints. 1081

- **9-point anchor:** End-to-end and logically 1082
tight; clearly specifies inputs/outputs; proac- 1083
tively addresses exceptions, edge cases, and 1084
includes fallback or contingency plans. 1085

(4) Final impact.

- **3-point anchor:** Fails to solve the task or 1087
introduces new bugs/safety issues; produces 1088
negative or misleading outcomes. 1089

- **6-point anchor:** Resolves the explicit prob- 1090
lem but results are mediocre or rough; “works” 1091
without notable quality, efficiency, or clarity. 1092

- **9-point anchor:** High-quality outcome that 1093
exceeds baseline expectations; efficient, clean, 1094
and well-structured; demonstrates strong pro- 1095
fessional judgment and execution. 1096

(5) Scalability / extensibility.

- **3-point anchor:** Highly coupled and hard- 1097
coded; brittle; one-off solution that is difficult 1098
to reuse or adapt; small changes require large 1099
rewrites. 1100
1101

1102 • **6-point anchor:** Moderately modular; ex- 1143
 1103 poses basic parameters/interfaces; can accom-
 1104 modate a reasonable range of requirement
 1105 changes.

1106 • **9-point anchor:** Highly decoupled and mod- 1145
 1107 ular; designed for reuse and iteration; easy
 1108 to migrate, extend, and evolve into a general
 1109 component.

1110 **(6) Alignment with task requirements.**

1111 • **3-point anchor:** Severely violates explicit 1147
 1112 constraints, uses prohibited tools/methods, or
 1113 is substantially off-topic.

1114 • **6-point anchor:** Satisfies the core require- 1150
 1115 ments but has minor deviations in secondary
 1116 constraints (e.g., format, length, or non- 1151
 1117 critical instructions).

1118 • **9-point anchor:** Strictly follows all explicit 1152
 1119 and implicit instructions; matches constraints
 1120 with near machine-level precision.

1121 **Self-consistency constraint.** We apply a judg- 1152
 1122 ing consistency rule: when a solution is rated as
 1123 low feasibility, its final impact should typically not
 1124 exceed the high-point anchor, to avoid rewarding
 1125 impractical but superficially impressive proposals.

1126 **D.4 Stability and Consistency Analysis**

1127 To assess the reliability of LLM-as-a-Judge eval- 1153
 1128 uation, we analyze **score consistency** and **rank-**
 1129 **ing stability** under repeated judging. Concretely,
 1130 we evaluate **7 methods** and repeat the judging **10**
 1131 **times** (by re-randomizing candidate order and re-
 1132 invoking the judges under the same unified prompt),
 1133 yielding $K = 10$ repeated outcomes per method.

1134 **Score consistency (ICC).** Score-level consis- 1154
 1135 tency measures whether the *absolute* scores as-
 1136 signed to the same method remain stable across
 1137 repeated judging. Let x_{ij} denote the scalar score of
 1138 method $i \in \{1, \dots, N\}$ in repeat $j \in \{1, \dots, K\}$,
 1139 where $N = 7$ and $K = 10$. Define means

1140
$$\bar{x}_{i\cdot} = \frac{1}{K} \sum_{j=1}^K x_{ij}, \quad \bar{x}_{\cdot j} = \frac{1}{N} \sum_{i=1}^N x_{ij},$$

1141
 1142
$$\bar{x}_{\cdot\cdot} = \frac{1}{NK} \sum_{i=1}^N \sum_{j=1}^K x_{ij}.$$

We compute the two-way ANOVA mean squares

1144
$$MS_R = \frac{K \sum_{i=1}^N (\bar{x}_{i\cdot} - \bar{x}_{\cdot\cdot})^2}{N - 1},$$

1145
 1146
$$MS_E = \frac{\sum_{i=1}^N \sum_{j=1}^K (x_{ij} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x}_{\cdot\cdot})^2}{(N - 1)(K - 1)}.$$

1147 We report the intraclass correlation coefficient
 1148 (ICC) for single-measure agreement:

1149
$$\text{ICC} = \frac{MS_R - MS_E}{MS_R + (K - 1)MS_E},$$

1150 where a higher ICC indicates more repeatable (less
 1151 noisy) scores across the 10 repeated judging runs.

1152 **Ranking stability (Kendall’s W).** Ranking sta-
 1153 bility measures whether the *relative ordering* of
 1154 methods is consistent across repeated judging. Let
 1155 r_{ij} be the rank of method $i \in \{1, \dots, M\}$ in re-
 1156 peat $j \in \{1, \dots, K\}$, with $M = 7$ and $K = 10$
 1157 (smaller rank is better). Define rank sums

1158
$$R_i = \sum_{j=1}^K r_{ij}, \quad \bar{R} = \frac{1}{M} \sum_{i=1}^M R_i,$$

1159 and dispersion

1160
$$S = \sum_{i=1}^M (R_i - \bar{R})^2.$$

1161 We compute Kendall’s coefficient of concordance

1162
$$W = \frac{12S}{K^2(M^3 - M)} \in [0, 1],$$

1163 where larger W indicates stronger agreement
 1164 among the 10 repeated rankings (i.e., higher rank-
 1165 ing stability).

1166 From Table 6, across all four judging LLMs,
 1167 score-level reliability is consistently high: the
 1168 **ICC(3,10)** values exceed **0.94** for most dimensions,
 1169 indicating that repeated judging produces highly re-
 1170 producible scores. The strongest consistency is ob-
 1171 served for Completeness, Final impact, and Scala-
 1172 bility (typically ≥ 0.97). Ranking-level stability is
 1173 also strong, with Kendall’s W in $[0.8407, 0.8736]$
 1174 across judges, implying that the relative ordering
 1175 of the 7 methods remains largely consistent over
 1176 the 10 repeated evaluations. Overall, these results
 1177 support that our LLM-as-a-Judge protocol is stable
 1178 in both absolute scoring and comparative ranking.

Table 6: Judge consistency and ranking stability in LLM-as-a-Judge evaluation.

Judging LLM	Grok-4-reason	Deepseek-V3.2	Gemini-3-pro	GPT-5.1-think
<i>Score consistency ICC (3,10)</i>				
Cost-effectiveness	0.9777	0.9798	0.9643	0.9494
Feasibility	0.9424	0.8720	0.9389	0.9588
Completeness	0.9811	0.9807	0.9729	0.9785
Final impact	0.9687	0.9856	0.9809	0.9744
Scalability	0.9807	0.9799	0.9761	0.9678
Alignment with task	0.8805	0.7882	0.9365	0.9386
<i>Ranking stability (Kendall's W)</i>				
Rank	0.8736	0.8407	0.8407	0.8593

E Description and statistical results of the CLM dataset

To evaluate long-horizon planning under complex, coupled constraints, we construct **CLM** (*Complex Long-horizon Multi-task Benchmark*), a curated benchmark of **50** high-difficulty task scenarios. CLM is designed around intrinsic interdisciplinarity: each task reflects a real-world techno-social engineering problem whose feasible solution must jointly integrate *physical*, *informational*, and *socio-economic* factors, rather than being solvable by single-domain knowledge alone.

CLM tasks are constructed following three principles:

- (1) each task involves multiple interdependent constraints that cannot be verified independently;
- (2) successful solutions require coordination across at least three roles or subtasks;
- (3) no single agent can trivially enumerate a complete solution in one step.

E.1 Taxonomy and Macro-level Distribution

CLM adopts a two-level taxonomy for systematic organization. At the top level, **MacroSystem** defines six pillars of modern techno-social systems; at the second level, **EcoSystem** refines each MacroSystem into application subdomains. This taxonomy is defined by the *primary application context* of each task and does not imply domain isolation—cross-domain coupling is a central feature throughout the dataset. The six MacroSystems and their counts are: *Energy & Enviro* (12/50, 24%), *Computing & Network* (10/50, 20%), *Smart Cities* (7/50, 14%), *Life Sciences* (7/50, 14%), *Social & Digital Finance* (7/50, 14%), and *Aerospace & Physics* (7/50, 14%).



Figure 5: Distribution of CLM Dataset

E.2 Within-task Cross-domain Coupling and Constraints

A defining characteristic of CLM is within-task cross-disciplinary coupling: most tasks require simultaneously invoking two or more distinct knowledge bases, where feasibility emerges from logical entanglement rather than simple knowledge aggregation. Representative coupling patterns include:

- (i) physics–perception–chemistry coupling (e.g., mapping nonlinear chemical mixtures to psychophysical perception thresholds, or closing real-time loops between thermodynamic constraints and learned inversion in industrial processes);
- (ii) law–logic–code coupling (e.g., translating international commercial law into deontic-logic code with formal verification);
- (iii) biology–cryptography–finance coupling (e.g., genomic pricing that jointly requires actuarial risk modeling, bioinformatic mutation analysis, and provable privacy/security constraints).

These patterns intentionally stress cross-domain

1235 coordination over isolated domain expertise.

1236 Each CLM task embeds multi-dimensional cou-
1237 pled constraints, with **3–4 hard requirements per**
1238 **task on average**. We group constraints into four
1239 major categories:

1240 (i) physical & temporal constraints (strict latency,
1241 dynamics, or physics limits),

1242 (ii) safety & robustness constraints (stability un-
1243 der failures, attacks, or uncertainty),

1244 (iii) economic & resource-efficiency constraints
1245 (cost, energy, or resource budgets),

1246 (iv) ethical & compliance constraints (fairness,
1247 privacy, and legal requirements).

1248 These constraints are intentionally coupled: sat-
1249 isfying only a subset typically yields solutions that
1250 are incomplete or logically infeasible, thereby re-
1251 quiring globally consistent, long-horizon reason-
1252 ing.

1253 **E.3 CLM Problem Examples and Output** 1254 **Prompt**

1255 CLM Problem examples and output prompts are
1256 shown in Figure 6 and Figure 7.

CLM Task Input Example: Capital Allocation under Coupled Constraints

Background.

Zhilian Cloud Services Co., Ltd. is a mid-sized enterprise SaaS provider with foreign investment background. The company focuses on enterprise-level cloud services and operates in a highly competitive cloud infrastructure market. Its core offering is basic cloud storage services, while it is actively expanding into emerging markets and AI-integrated product lines. Facing intense competition from major cloud incumbents, the company urgently seeks differentiated growth strategies and brand positioning. The company plans to allocate funds from an overseas financing round of RMB 20 million. Among these, RMB 12 million is designated for rapid validation of three business lines to explore growth boundaries. The founder explicitly requires that at least RMB 3 million must be reserved for the core business to ensure service-level agreement (SLA) compliance and operational security for existing customers.

Business Lines.

- **A. Core Business: Enterprise Cloud Storage Subscription (Low Risk)**
Provides standardized cloud storage and backup services to small and medium-sized enterprises on an annual subscription basis. Customer churn rate is below 5%, and this line generates stable cash flow and serves as the primary profit source.
- **B. Market Expansion: Latin American SME Penetration and High-end Industry Solutions (Medium Risk)**
Investment focuses on building local compliance teams, participating in industry summits, and delivering customized solutions. Target verticals include manufacturing and retail sectors.
- **C. R&D / Breakthrough Project: Generative AI Customer Service Engine (High Risk, High Return)**
Aims to develop a multilingual generative AI customer service model for productization and large-enterprise procurement. Funding is used for data acquisition, model iteration, and global patent deployment.

Task.

Given the above context, determine how the company should allocate capital across the three business lines to address current challenges while enabling sustainable future growth. The solution must respect the minimum funding constraint for the core business and balance risk, return, and strategic optionality.

Figure 6: CLM Task Input Example

System Output Prompt (Assistant)

As a solution expert, you must:

1. Provide a complete, structured solution in Markdown.
2. Your output must include the following sections:
 - **Problem Analysis:** your core understanding of the user's needs
 - **Solution:** detailed, actionable step-by-step instructions
 - **Key Points:** critical considerations and core logic
 - **Expected Results:** a clear description of the deliverables
3. If code examples are needed, use a standard Markdown code block in the following format:

```
```python
code comments
implementation
```

4. Do not output any explanatory preface or any concluding remarks.
5. Do not add any extra characters (e.g., # or / ) before Markdown markers.
6. The content must be directly usable; avoid placeholders and ellipses.
7. Use clear, concise professional language.
8. Output the complete Markdown solution only, with no additional text.

Figure 7: System Output Prompt