

Know-the-Ropes: Algorithmic Blueprints for Reliable LLM Multi-Agent System Design

Anonymous ACL submission

Abstract

Single-agent LLMs face finite context and role overload, while unstructured multi-agent designs can introduce ambiguous roles and coordination overhead. We therefore introduce Know-The-Ropes (KtR), a practical methodology for projecting algorithmic priors and heuristics into typed, controller-mediated multi-agent blueprints for decomposable tasks. KtR follows a multi-step process—identify bottlenecks, refine decomposition, apply minimal augmentation (chain-of-thought, self-check, or light fine-tuning), and verify via contracts. In two case studies, including Knapsack (3–8 items) and Task Assignment (6–15 jobs), we find that KtR by low-effort LLMs can show notable end-to-end accuracy gains over single-agent zero-shot baselines. With three GPT-4o-mini agents, accuracy on size-5 Knapsack instances rises from 3% to 95% after addressing a single bottleneck agent. With six o3-mini agents, Task Assignment reaches 100% up to size 10 and $\geq 84\%$ on sizes 13–15, versus $\leq 11\%$ zero-shot. These results indicate benefits in our controlled setting; KtR complements scaling and prompt/program-of-thought techniques in building a reliable multi-agent system. An anonymous code base is available at [this anonymous link](#)

1 Introduction

Large language model (LLM) agents typically achieve strong performance in the domains for which they are trained and optimized (Thirunavukarasu et al., 2023; Kasneci et al., 2023; Wu et al., 2023b), yet their effectiveness degrades outside those boundaries. Finite context windows limit long-document reasoning, and no single agent can simultaneously address mathematics, coding, and planning problems (Liu et al., 2023; Gulati et al., 2024; Wu et al., 2025). Persistent challenges, including hallucinations, topical drift, and domain-specific failures, further constrain their reliability (Zhang et al., 2024; Xu et al.,

2024). A natural direction is the division of labor by decomposing tasks across specialized agents that coordinate to produce a joint solution. Current frameworks such as Mixture-of-Agents suggest that, when properly orchestrated, a team of agents can outperform even its strongest individual member (Wang et al., 2024a).

Yet significant challenges remain. First, each task still demands a carefully crafted prompt, often requiring substantial manual effort (Li et al., 2025b; Cui et al., 2025). While some multi-agent frameworks report encouraging results, once evaluation leakage and prompt overfitting are controlled, the apparent gains of naïve agent swarms collapse to single digits, and can even turn negative when tasks require more rounds of coordination (Pan et al., 2025; Zhu et al., 2025). Post-hoc analyses consistently reveal recurrent failure modalities: ill-posed task decompositions propagate ambiguity, imprecisely defined roles lead to redundancy or coverage gaps, and verification mechanisms are either predicated on brittle heuristics or become computationally prohibitive (Cao et al., 2025; Wang et al., 2024b). Furthermore, latency and cost tend to scale super-linearly with each round of interaction (Ye, 2025; Shu et al., 2024). These findings suggest that simply adding more “brains” does not guarantee progress. Robust and scalable improvements of multi-agent system (MAS) design demands a principled, systems-engineering approach—the gap our work aims to close.

To address this critical need, we introduce Know-The-Ropes (KtR), a strategy that reframes MAS design as structured algorithm engineering. At its core, KtR employs a hierarchical task decomposition, recursively partitioning problems into sub-tasks until each primitive operation is solvable by a base model, potentially augmented with minimal augmentation (e.g., self-check loops, Chain-of-Thought (CoT) (Wei et al., 2022), fine-tuning). Inter-agent communication is not conversational

085 but is instead mediated by a code-based controller
086 that enforces explicit, typed input/output contracts.
087 This architectural choice ensures modularity and
088 predictable information flow, preventing common
089 pathologies such as context bloat and state over-
090 writes. This design philosophy turns MAS con-
091 struction into traditional algorithm engineering:
092 identify bottlenecks, refine the decomposition, and
093 apply the most cost-effective augmentation to un-
094 derperforming agents, mirroring the optimization
095 of classical computational graphs.

096 Our empirical study includes the following ex-
097 periments: **(i) Knapsack Problem.** Starting with
098 GPT-4o-mini, zero-shot accuracy on 3–8-item in-
099 stances ranges from 60% to 0%. A naïve task-level
100 fine-tune shows limited improvement. Under our
101 KtR three-agent blueprint, accuracy improves to
102 95%–70% in our controlled setting after fine-tuning
103 just the “trimmer” sub-task on 1200 worked exam-
104 ples. **(ii) Task-Assignment Problem (scalability
105 case study).** With o3-mini we build a six-agent
106 blueprint for problems of size 6–15. Decomposing
107 a single weak agent into two finer leaves yields
108 leaf accuracies of 100% and 97%, with the overall
109 system achieving $\geq 84\%$ accuracy across all sizes
110 in our evaluation protocol. Based on these findings,
111 our contributions lie in three aspects.

- 112 • **A Strategy for MAS Design.** We present
113 KtR, a strategy for translating algorithmic pri-
114 ors into multi-agent blueprints with typed in-
115 terfaces and local verification, applicable to
116 decomposable tasks with known domain struc-
117 ture.
- 118 • **Empirical Validation.** Two illustrative
119 demonstrations (Knapsack and Task-
120 Assignment problems) indicate that targeted
121 decomposition can improve end-to-end
122 accuracy versus single-agent baselines in our
123 settings, using modest models and minimal
124 augmentation.

125 2 Related Work

126 Multi-Agent Systems (MAS) have been widely
127 employed to enhance the capabilities of LLMs to
128 tackle complex tasks (Qiu et al., 2024; Yan et al.,
129 2024; Ma et al., 2024; Lin et al., 2024; Hua et al.,
130 2023; Yu et al., 2024). This is because MAS typi-
131 cally distribute tasks across agents that collaborate
132 to achieve a common goal, thereby improving both
133 efficiency and adaptability. Recent frameworks
134 like CAMEL (Li et al., 2023) enable role-based
135 cooperative dialogues by assigning agents distinct

136 personas, while AutoGen (Wu et al., 2023a) and
137 MetaGPT (Hong et al., 2023) orchestrate multi-
138 role agent teams through structured conversation
139 loops and predefined workflows. In math opti-
140 mization, OR-LLM-Agent can translate natural-
141 language problem descriptions into formal Gurobi
142 models—achieving an 85% correct-solution rate
143 on real-world benchmarks (Zhang and Luo, 2025).

144 However, studies show that simply scaling up to
145 LLM-based MAS often yields only marginal gains
146 over single-agent baselines (Pan et al., 2025). LLM
147 agents still struggle with context management and
148 consistency, meaning that elaborate multi-agent
149 prompts can fail to realize the intended collabora-
150 tion (Bo et al., 2024). A recent systematic audit of
151 popular MAS frameworks has identified 14 distinct
152 failure modes (Cemri et al., 2025), which can be
153 grouped into three categories, including flawed de-
154 sign (e.g., ambiguous role definition), inter-agent
155 misalignment (e.g., communication failures), and
156 quality control (e.g., no reliable check mechanism).

157 To address these challenges, researchers have
158 proposed multiple strategies to make LLM-based
159 MAS more reliable (Zhu et al., 2025; Tran et al.,
160 2025). A key strategy is improving the agent in-
161 teraction structure (Zhu et al., 2025). For example,
162 the AgentDropout framework proposes a dynamic
163 agent-pruning strategy, which seeks to discard less
164 critical actors during training (Wang et al., 2025).
165 Another effective strategy is incorporating feed-
166 back and verification loops (Hong et al., 2023). A
167 recent study shows that frameworks with role spe-
168 cialization and iterative feedback mechanisms can
169 outperform those without these features (Anony-
170 mous, 2025). In addition, systematic evaluations
171 suggest that the communication topology matters:
172 a well-designed protocol between agents can signif-
173 icantly improve collective performance on complex
174 tasks (Zhu et al., 2025).

175 While existing multi-agent frameworks and
176 strategies have demonstrated notable progress (Li
177 et al., 2025a; Hong et al., 2023; Zhu et al., 2025),
178 they still face challenges in dynamic role real-
179 location and efficient inter-agent communication.
180 These limitations become especially pronounced
181 when addressing complex tasks, such as NP-hard
182 optimization problems. To bridge this gap, our
183 work proposes a heuristic strategy that embeds
184 domain-specific rules and algorithms directly into
185 agent coordination. This approach enables on-the-
186 fly role adaptation and task decomposition, particu-
187 larly in math optimization contexts where conven-

188 tional MAS frameworks often struggle.

189 3 Methodology-KtR Framework Design

190 We propose the strategy "Know the Ropes" (KtR),
191 which offers a structured methodology for de-
192 signing specialized MAS leveraging LLMs. This
193 heuristic focuses on translating known, effective
194 procedures or algorithms into a coherent multi-
195 agent architecture. As presented in Figure 1, the
196 core idea is to decompose a complex overall task
197 into its fundamental computational stages. Each
198 stage is then mapped to a well-formulated sub-task,
199 designed to be tractable for an individual agent.
200 These specialized agents are subsequently orches-
201 trated to mirror the data/control flow of the original
202 procedure, which can effectively embed problem-
203 solving logic into the MAS. The following defini-
204 tions formalize the KtR components.

205 This approach is grounded in the No-Free-Lunch
206 (NFL) theorem (Wolpert and Macready, 1997;
207 Wolpert, 2021) (see Appendix C), which states that
208 no algorithm performs universally better than oth-
209 ers across all problem distributions. Rather than
210 seeking a universal multi-agent orchestration strat-
211 egy, KtR operationalizes the NFL insight by inject-
212 ing domain-specific inductive bias through algo-
213 rithmic blueprints. In particular, domain-specific
214 heuristics is used to guide us to properly decom-
215 pose the task, and by decomposing tasks according
216 to their inherent algorithmic structure, we concen-
217 trate the system’s “learning budget” on the specific
218 problem distribution at hand, trading generality for
219 reliability within the target domain.

220 To formulate our idea in some generality, we will
221 introduce and use several bolded terms below, with
222 only a brief explanation. More formal definitions
223 for these terms can be found in Appendix.

224 (i). We start with a **well-formulated** task, which
225 we denote by a triple (I, O, R) where I and O
226 represents input and output of the task, and R out-
227 lines a clear rule how we obtain desired output
228 from the input, we want to produce an MAS de-
229 sign, which we call a **workflow blueprint**, and
230 denote by $\mathcal{B} = (\mathcal{T}, \mathcal{P})$. Here \mathcal{T} is a collection of
231 (well-formulated) tasks and \mathcal{P} is the protocol (the
232 hard-coded workflow) to concatenates the I/O of
233 consecutive tasks to solve the original one.

234 (ii). The method to obtain a practical workflow
235 blueprint is via **task decompositions** that is guided
236 by the domain-specific heuristics. This is a recur-
237 sive process: we test if each current task in \mathcal{T} is

238 **tractable**, i.e., if a single agent, with necessary
239 augmentations such as prompt engineering or fine-
240 tuning can handle the task with high accuracy. The
241 decomposition is then guided by two aspects:

- 242 • The agent-wise tests which identifies the bot-
243 tleneck, indicating which task to be decomposed.
- 244 • The domain-specific heuristics guide us how
245 to decompose.

246 We repeat the process until all tasks become
247 tractable. We call this recursive testing-
248 decomposing process a **tractable hierarchy**.

249 (iii). After acquiring a tractable hierarchy, we
250 then map the final blueprint into a practical MAS
251 system: each (tractable) tasks is now handled by
252 an agent, and the protocol correspond to how the
253 system controller orchestrates the agents.

254 This three-step procedure, algorithmic blueprint,
255 tractable hierarchy construction, and system in-
256 stantiation, provides a principled pathway from
257 a complex task to a deployable MAS solution with
258 correctness hinges on model capabilities that have
259 been explicitly validated. To demonstrate the prac-
260 tical application and efficacy of the KtR framework,
261 we investigate two case studies. In each case, we
262 use a well-understood algorithm to decompose the
263 complex problem into an M-tractable hierarchy and
264 instantiated MAS.

Algorithm 1 KtR Framework Pseudo-code

```
1: procedure KTR( $T, \mathcal{M}$ )
2:    $B \leftarrow \text{CREATEBLUEPRINT}(\{T\}, \text{trivial\_protocol})$  #
   Start with the top-level task
3:   while exists  $U \in B.\text{tasks}$  and
    $\neg \text{MTRACTABLE}(U, \mathcal{M})$  do
4:      $U^* \leftarrow \text{CHOOSETASKTODECOMPOSE}(U)$  #
     Select a non-tractable task
5:      $B_{\text{sub}} \leftarrow \text{DESIGNSUBBLUEPRINT}(U^*)$  # Define
     its sub-tasks and protocol
6:      $B \leftarrow \text{EMBEDSUBBLUEPRINT}(B, U^*, B_{\text{sub}})$  #
     Replace the task with its sub-blueprint
7:      $\text{ASSERTINTERFACEPRESERVED}(B)$  # Ensure
     communications are valid
8:   end while
9:    $\text{MAS} \leftarrow \text{INSTANTIATESYSTEM}()$  # Begin building
   the MAS
10:  for all  $V \in B.\text{tasks}$  do
11:     $\text{aug} \leftarrow \text{SELECTAUGMENTATIONS}(V, \mathcal{M})$  #
    Select a cost-effective augmentation
12:     $\text{agent} \leftarrow \text{CREATEAGENT}(V, \mathcal{M}, \text{aug})$  # Create
    a specialized agent based on the definition
13:     $\text{MAS.AddAgent}(\text{agent})$  # Add the specialized
    agent to the system
14:  end for
15:   $\text{IMPLEMENTPROTOCOL}(\text{MAS}, B.\text{protocol})$  # Wire
    up agents based on the blueprint
16:  return  $\text{MAS}$  # Return the final multi-agent system
17: end procedure
```

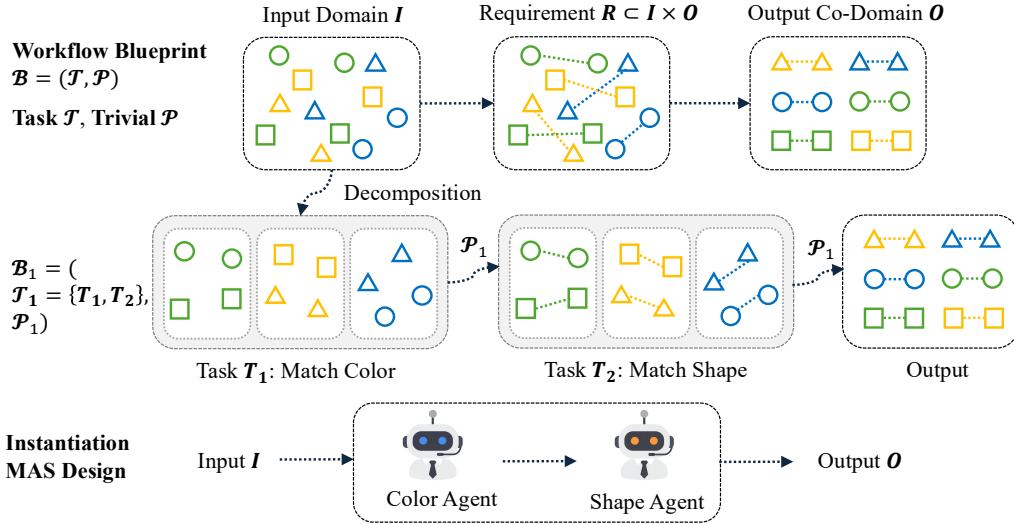


Figure 1: Illustration of the KtR strategy: heuristic, prior-guided decomposition of a complex task into sub-tasks, each instantiated as a coordinated LLM agent within a multi-agent architecture.

4 Experiment Design

4.1 Proof-of-Concept: 0/1 Knapsack Problem (KSP)

To furnish a clear proof-of-concept for KtR, we start with the classical NP-hard Knapsack Problem (KSP), a staple in resource allocation, logistics, and investment planning. Using the lightweight, general-purpose GPT-4o-mini as the MAS backbone, we establish a modest baseline that allows us to highlight how KtR MAS choreography amplifies a small model’s capability well beyond its limits.

Problem Formulation. For KSP, the input is a tuple (\vec{w}, \vec{v}, W) where \vec{w} and \vec{v} are two N -dimensional vectors representing the weight and value of n items, and W is the weight capacity. Then the objective of the Knapsack problem can be formulated as finding the optimal value:

$$Z = \max\{\vec{x} \cdot \vec{v} \mid \vec{x} \in \{0, 1\}^N, \vec{x} \cdot \vec{w} \leq W\}.$$

Here $\vec{x} \in \{0, 1\}^N$ is the state vector, representing whether an item is chosen or not. This problem, where each item can either be fully included or not at all, is commonly known as the 0/1 KSP.

Problem Solution A classic approach to the Knapsack Problem iteratively enumerates all feasible states—a dynamic programming strategy (Bellman, 1957). Formulated in the form of mathematical induction, the initial state is $S_0 = \{(0, 0)\}$, and we add items in inductively, with capacity being aware: for each k , assuming that S_{k-1} has been obtained, then we form:

$$S_{add} = \{(w + w_k, v + v_k) \text{ for all } (w, v) \in S_{k-1}\}.$$

Then we trim by the capacity $S_{trimmed} = \{(w, v) \in S_{add} \mid w \leq W\}$ and union the two set of states to create $S_k: S_k = S_{k-1} \cup S_{trimmed}$. After running through all items and obtaining the set S_N , we pick the element in S_N with maximal value, as the solution to the KSP. Specific details of this method is presented in Appendix D.

KtR MAS Design. Following the “Know the Ropes” heuristic, the iterative dynamic programming solution for the KSP is decomposed into tasks for three specialized agents, including: (i) **Worker Agent:** Computing the set S_{add} from S_{k-1} and the k -th item (w_k, v_k) . (ii) **Trimmer Agent:** Obtain $S_{trimmed}$ from S_{add} and the capacity W . (iii) **Reporter Agent:** Find the element with maximal value within the final state set S_N . (iv) **System Controller:** The controller orchestrates the overall process. After initialization, it controls the loop on k . For each k , it sends S_{k-1} and (w_k, v_k) to Worker Agent, and then send the result plus W to the Trimmer Agent. The Controller then takes the union to obtain S_k . Once all items are processed, the Controller invokes the Reporter Agent for final result. Specific prompts for these agent design are attached in Appendix G.

4.2 Proof of scalability—Task Assignment Problem (TAP)

Building on the previous section where KtR already stretched the capabilities of the compact GPT-4o-mini on the Knapsack baseline, we now test the framework’s scalability. We upgrade the backbone to the larger o3-mini and tackle the more demand-

ing Task-Assignment Problem (TAP), demonstrating that the framework’s performance rises in lock-step with the underlying model’s capacity.

Problem Formulation. TAP seeks to optimally assign a set of N workers to N tasks, where each potential assignment incurs a specific cost to minimize the total cost. The input is an $N \times N$ matrix C representing the cost. Then the objective of the TAP is to find the optimized value:

$$Z = \max_{\sigma \in \mathfrak{S}_N} \sum_{i=1}^n C_{i\sigma(i)}.$$

Here \mathfrak{S}_N represents the set of all permutations of N elements, representing arrangements that assign different tasks to different agents.

Problem Solution. The typical solution for TAP is using the Hungarian algorithm (Kuhn, 1955), which provides a polynomial-time method to find the objective value Z . We summarize the algorithm as follows. **Step 1. Row Reduction.** For each row, we subtract each entry with the minimal value of all entries. This create a matrix C' . **Step 2. Column Reduction.** Same operation for each column to obtain a new matrix C'' . **Step 3. Zero Covering.** We then find the smallest collection \mathcal{C} of rows and columns to cover all zeroes. Let L be the size. If $L = N$, then we skip Step 4. **Step 4. Matrix Improvement.** If $L < N$, we find the minimal value m of all entries that are not covered by the collection \mathcal{C} , and then subtract m from all uncovered entries and add m to all covered entries. Let C''' be the resulting matrix. **Step 5. Assignment Identification.** If $L = N$, then we attempt to find a collection of zeroes in C'' or C''' in which no two are on the same row or column. The position of the zeroes represents the optimal task assignment. Specific details of this method is presented in Appendix D.

KtR MAS Design. We apply the KtR methodology to create the MAS. As explained in Section 5.2, based on test results of the agentic tasks and heuristic argument, we further decompose step 3 into two agents. For better presentation, let N be the size of the original TAP problems, i.e., the number of tasks and resources in the problem. (i) **Row Reducer:** Realizes Step 1 and obtain the reduced matrix C' . (ii) **Column Reducer:** Realizes step 2 and obtain C'' . (iii) **Matcher:** Find a maximal collection of zeroes in the reduced matrix C'' in which no two are on the same row or column. Let L be the size of the collection. (iv) **Painter:** Find a minimal collection of rows and columns covering all zeroes when $L < n$. (v) **Normalizer:** Creates more ze-

roes outside of the selected rows and columns to get C''' . (iv) **Reporter:** Report the final answer when $L = N$. (vii) The **System Controller:** Arranges task for Row Reducer and Column Reducer linearly, then controls a loop: Matcher finds a set of zeroes and then Controller checks the size to determine whether to break the loop or not. In the loop, Painter is then called in to find the collection of lines and Normalizer follows to create more zeroes. Outside the loop, the Reporter is called to deduce the final answer. Specific prompts for these agent design are attached in Appendix G.

5 Experiment Result

Our experimental protocol unfolds in two stages. First, we run a uniform benchmark across a suite of baseline models—including several GPT and Llama variants—to fix a reference point for each task. The second stage then splits by objective: For KSP we deliver a proof of concept, while for TAP we provide a proof of scalability. For ground truth, we use python code to randomly generate problems, and then use the Google OR-Tools (Perron and Furnon, 2022) as in Appendix F to generate solutions to compare with.

5.1 Experiment Result for KSP

Baseline LLM Performance Figure 2 shows the baseline LLM performance across multiple difficulty levels. The accuracy across difficulty levels (from 3 to 8 items) in the KSP scenario reveals substantial performance variation among the tested LLMs. Among those, the GPT-o3-mini, as a reasoning model, consistently demonstrates superior accuracy. Model GPT-4.1 outperforming its smaller counterparts, namely GPT-4.1-mini and its primer GPT-4o-mini. Other LLMs, including Claude-haiku, Llama, and Qwen series also show performance degradation, with higher variability particularly evident at greater difficulty levels. Meanwhile, the performance of final KtR MAS is also drawn in Figure 2. The comparison shows that KtR substantially boosts performance, validating its effectiveness.

KtR MAS Performance Based on Figure 2, GPT-4o-mini exhibits a pronounced performance decline beginning at instances of 4 items, underscoring its limited scalability to more complex scenarios; therefore, we select it as the backbone for our KtR framework design. Figure 3 further illustrates the resulting KtR MAS along with the

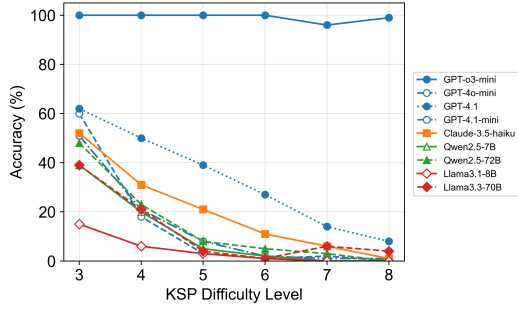


Figure 2: KSP baseline performance from single LLMs as well as the KtR MAS.

experimental outcomes based on our strategy.

(i) Single LLM performance. We establish two baseline performances for GPT-4o-mini acting as a single agent to solve the KSP. First, the zero-shot GPT-4o-mini is directly prompted with KSP instances. As Figure 3B shows, its accuracy decreases from 60% for 3 items to 0% (8 items). Second, we evaluate a fine-tuned GPT-4o-mini (standalone). Figure 3C indicates that fine-tuning offers some improvement over the zero-shot, but still as low as 3% for 8-item KSP.

(ii) Standard MAS performance. Following our KtR heuristic, we map the algorithm for KSP into a MAS design, illustrated in Figure 3F. Initially, each agent is driven by the standard, non-fine-tuned GPT-4o-mini. The performance of this standard MAS is presented in Figure 3F. Its performance decreases from 18% for 3 items to 4% for 8 items. This initial result implies that without augmenting the agents’ abilities, the MAS does not effectively handle the task. We profile each agent in isolation (Figure 3E) and uncover a single choke point: Trimmer. Its accuracy collapses as the feasible-state set S_k (cf. Section D.2) grows—54 % for 1–8 states, 24 % for 9–16, 7 % for 17–24, and just 5 % for 25–32. Because the algorithm loops once per state, even small per-iteration errors compound, and this cascading inaccuracy ultimately sinks the entire run.

(iii) Augmented MAS performance. To eliminate the bottleneck, we fine-tune the Trimmer’s GPT-4o-mini backbone (Figure 3G, highlighted as ‘Augmented Trimmer’). Accuracy leapt to 95 % for 1–8 feasible states, 89 % for 9–16, 81 % for 17–24, and 67 % for 25–32 (Figure 3H). Adding a lightweight self-check—prompting the model to audit its own answer—preserved or marginally improved these gains. Replacing the bottleneck with the fine-tuned Trimmer lifts end-to-end KSP accuracy to near-saturation across sizes (Figure 3I):

95 % for 3-item instances, 90 % for 4, 95 % for 5, 85 % for 6, 76 % for 7, and 70 % for 8. A single targeted upgrade thus turns KtR into a consistently high-performing solver as the problem scales.

5.2 Experiment result on TAP

Baseline LLM Performance Figure 4 illustrates the baseline performance of multiple LLMs on the TAP task across multiple difficulty levels (from 3 to 8 tasks). The results reveal marked differences in model capabilities. The only reasoning model, GPT-4o-mini, consistently outperforms all others, exhibiting strong accuracy at lower difficulty levels, though its performance declines as task complexity increases. In contrast, GPT-4.1 demonstrates moderate but stable accuracy across all difficulty levels, surpassing its mini-sized counterparts. Other models, including Claude-3.5-Haiku, Qwen2.5, and Llama-3 variants, show intermediate performance with variability.

We observe that single-agent models (e.g., GPT-3-mini, GPT-4-mini, GPT-4.1) drop to 30-50% accuracy at TAP levels 7-8, while KtR MAS maintains steady performance near 100%, even surpassing reasoning models, demonstrating its exceptional robustness and generalization capabilities.

KtR MAS Performance Based on Figure 4, GPT-o3-mini consistently outperforms other LLMs across all evaluated tasks, making it our choice for subsequent experiments. Our goal is to assess the scalability of our proposed strategy and investigate how its performance evolves as task difficulty increases. Figure 5 illustrates the MAS design and corresponding experimental outcomes obtained using our heuristic-based approach.

(i) Single LLM performance. Again, we evaluate the baseline performance of using o3-mini as a single agent. The o3-mini model achieves a relatively high performance (83%) but decays quickly as in Figure 5B: 37% on problems of size 9, 21% on problems of size 12 and finally is reduced to 3% for problems of size 15.

(ii) Standard MAS performance and further decomposition. Guided by the Hungarian algorithm (Kuhn, 1955), our first KtR blueprint maps each step to a single agent; Step 3 from Section 4.2 relied on a lone Cover Seeker rather than the later ‘Matcher + Painter’ pair. This baseline already scored 98 % (size 6), 88 % (size 9), 78 % (size 12), and 56 % (size 15), validating the approach.

We then stress-test each agent on two bands,

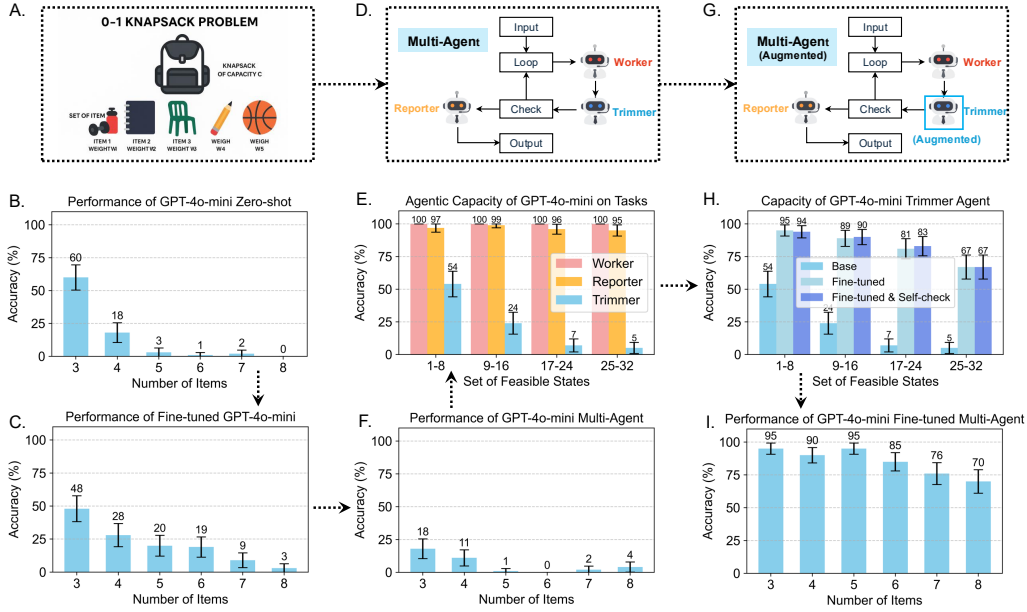


Figure 3: KSP evaluation of the KtR strategy. **B**: Zero-shot accuracy of the baseline model. **C**: Zero-shot accuracy after a light, task-specific fine-tune of the same model. **D & G**: Blueprints of the MAS without (**D**) and with (**G**) augmentations. **E**: Per-agent accuracies before augmentation, revealing the system’s bottleneck. **H**: Boost delivered by two targeted augmentations—task-level fine-tuning and self-check prompting—applied to the bottleneck agent. **F & I**: Corresponding test accuracies for the two blueprints.

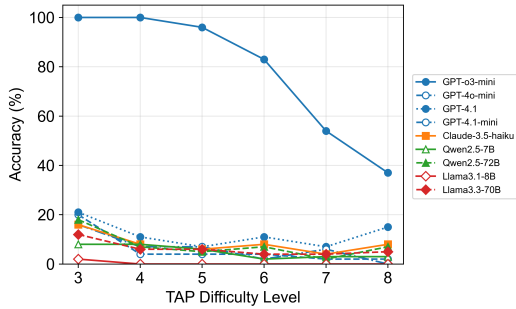


Figure 4: TAP baseline performance from single LLMs as well as the KtR multi-agent system.

with matrix sizes 6-10 and 11-15, to pinpoint weaknesses. One-shot agents are flawless: Row Reducer and Reporter reach 100% on both bands, and Column Reducer hit 100% / 92%. Normalizer holds 99% / 94%, but Cover Seeker falls to 97% / 84%. Because Zero Seeker operates inside the main loop, its errors accumulate, making it the clear bottleneck for larger TAP instances.

We then perform a further decomposition of Step 3 in Section D.2 in a two-step process: Step 3.1. Finding a **maximal** collection of zero-entries, such that no two share a same row or column; Step 3.2. Finding a **minimal** collection of rows and columns covering all zero-entries. We believe this decomposition is helpful due to the following reasons. First, by a mathematical argument, the size of collections from sub-tasks 3.1 and 3.2 match. Second,

a heuristic argument indicating that knowing the maximal collection of zeroes simplifies the task to find minimal collection of rows and columns. Last, the original Step 5 can then simply use the positions of the zeroes from Step 3.1, once optimal check passes. Note, this also explains why we prefer a further decomposition rather than fine-tuning the original agent, as a further decomposition improves the system flow as well. Empirical pays off, as shown in Figure 5I, Matcher reaches 100% accuracy on both difficulty bands, while Painter climbs to 98% and 97%—a sharp jump from the original Cover Seeker’s 97% and 84%.

(iii) Augmented MAS performance. Leveraging the refined decomposition, we deploy a six-agent system (Figure 5H) that solves size 6–10 instances almost flawlessly: almost 100% accuracy versus 83 – 27% for o3-mini zero-shot. It sustains high performance on size-11–15 tasks (95%, 97%, 90%, 93%, 84%); even the dip at size 15 far exceeds the 3% zero-shot baseline, highlighting the substantial capacity gain of our MAS.

(iv) Token and time usage In order to further compare the token and time usage of the multi-agent system design with the zero-shot, we have the following table on size 6, 9, and 12 game instances.

From Table 1, we can see that as the size (i.e. complexity) of the problem increases, the ratio of average costs between MAS and zero-shot grows.

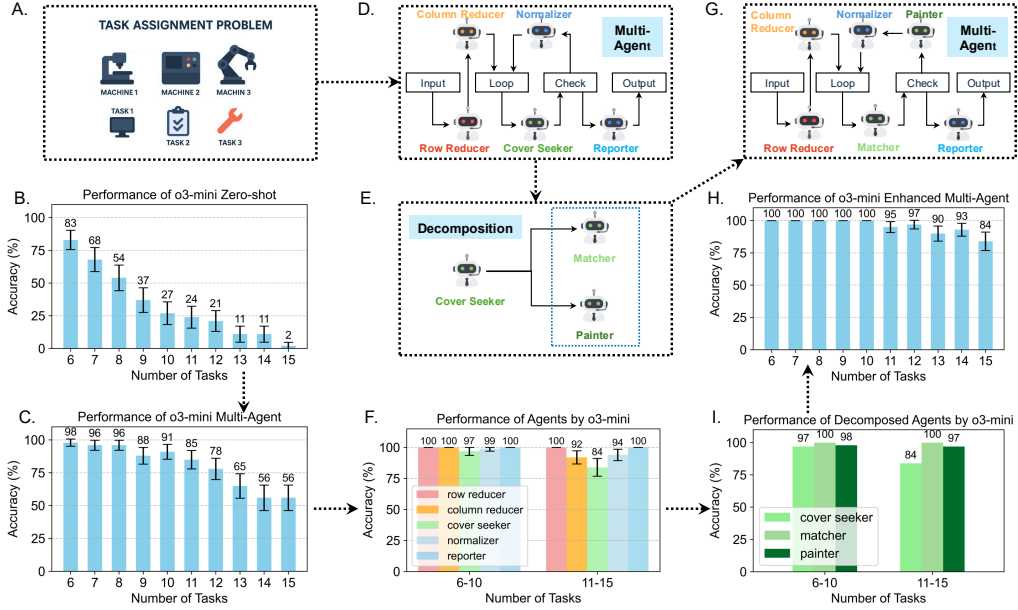


Figure 5: TAP evaluation of the KtR strategy. **B**: Zero-shot accuracy of the baseline model. **D**: Initial blueprint derived from the Hungarian algorithm; its end-to-end accuracy is shown in **C**. **F**: Per-agent accuracies within this blueprint, prompting the finer decomposition outlined in **E**. **I**: Side-by-side comparison of per-agent accuracies before and after decomposition. **G**: Final, decomposed blueprint, with overall accuracy presented in **H**.

In particular, for token usage, the input token ratio is around 10, while the output token cost ratio grows from 2.5 to 4. This indicates the trade between accuracy and consumptions.

Size	Solver	I. Token	O. Token	Time
6	Zero-shot	359.68	3338.40	21.61
	System	2392.09	8260.92	63.47
9	Zero-shot	494.78	6318.86	41.56
	System	4439.40	18618.88	132.05
12	Zero-shot	683.77	7906.11	48.40
	System	7290.43	32712.94	216.94

Table 1: Performance Statistics by Problem Size and Solver Type. "I. Token" represents the average input token cost, and "O. Token" represents the average output token cost. Time is the average time consumption per case.

6 Discussion and Conclusions

We present KtR, an engineering framework that turns algorithmic knowledge into reliable MAS via typed, verifiable blueprints. Our approach targets common failure modes in current MAS by enforcing structured decompositions, JSON-schema contracts, and local verification. Across evaluations, KtR shows effectiveness along complementary axes: KSP provides proof-of-concept by showing how a general-purpose small model (GPT-4o-mini) can handle complex reasoning tasks through structured decomposition, achieving 95% accuracy versus 3% zero-shot on our benchmark; TAP offers

proof-of-scalability by leveraging the stronger reasoning capabilities of o3-mini on more demanding instances, reaching 100% accuracy up to size 10. To assess the generalizability, we include a third case study on a real-world language task (e.g., academic writing proofread), detailed in Appendix E.

KtR’s principled “**identify** → **improve** → **verify**” methodology indicates that disciplined decomposition plus targeted augmentation can substantially improve underperforming models on these tasks where both single-agent and unstructured MAS approaches struggle. The framework’s effectiveness follows from domain-aligned inductive bias rather than prompt tinkering, broadly consistent with insights suggested by the No-Free-Lunch theorem. Our systematic process—bottleneck identification, task tractability assessment, and minimal augmentation—helps turn modest models into more reliable collaborators on decomposable problems without requiring ever-larger monoliths.

While our current validation spans structured optimization and language tasks, future work will extend to address real-world complexities including noisy inputs and automated bottleneck detection.

Next, KtR positions algorithmic blueprints as a reliability-first complement to model scaling and prompt design, offering a systematic pathway from complex problems to deployable MAS solutions for decomposable tasks.

7 Limitations

Despite demonstrating sizable performance gains, our study has several limitations that should guide future work.

Narrow task scope. We evaluate Know-The-Ropes (KtR) on two canonical optimization problems (Knapsack and Task-Assignment). While chosen to illustrate proof-of-concept and proof-of-scalability, these tasks have well-structured objective functions and small action spaces. Though we further provide a language task performance as a supplement in Appendix E, the adaptability our our KtR design in full generality remain untested.

Synthetic data & idealized inputs. All problem instances are randomly generated and fully specified. Real-world inputs (noisy, partially observed, or adversarial) could degrade both decomposition quality and agent reliability.

Cost and latency trade-offs. We were only able to provide a partial analysis on the cost and latency trade-offs. Although this is already informative, a full token and time consumption record may be preferred in future works.

Bottleneck identification heuristic. We locate bottleneck subtasks via held-out accuracy screens; this assumes the availability of inexpensive ground-truth labels. Automated bottleneck detection without labels is an open problem.

8 Ethical considerations

This paper proposes Know-the-Ropes (KtR), a methodology for projecting algorithmic prior knowledge into typed, controller-mediated multi-agent system (MAS) blueprints, and evaluates it on synthetic combinatorial optimization tasks: the Knapsack Subset Problem (KSP) and the Task Assignment Problem (TAP). Our experiments do not involve human subjects, user data, or personally identifiable information. Problem instances are programmatically generated and ground-truth solutions are computed with Google OR-Tools (Apache 2.0), as described in Section F. Experiments use hosted large-language-model APIs (specifically GPT-4o-mini and o3-mini) for inference; we do not train or distribute any model weights. The scope of our study is limited to benign, synthetic tasks; we do not deploy agents in real-world settings. We are not aware of conflicts of interest related to this submission. An anonymous code repository is provided solely to facilitate inspection and reproduction and requires users to comply with

the terms of the referenced third-party tooling.

References

- Anonymous. 2025. [Code in harmony: Evaluating multi-agent frameworks](#). In *Submitted to CS598 LLM Agent 2025 Workshop*. Under review.
- Richard Bellman. 1957. *Dynamic Programming*. Princeton University Press, Princeton, NJ, USA.
- Xiaoho Bo, Zeyu Zhang, Quanyu Dai, Xueyang Feng, Lei Wang, Rui Li, Xu Chen, and Ji-Rong Wen. 2024. Reflective multi-agent collaboration based on large language models. *Advances in Neural Information Processing Systems*, 37:138595–138631.
- Pengfei Cao, Tianyi Men, Wencan Liu, Jingwen Zhang, Xuzhao Li, Xixun Lin, Dianbo Sui, Yanan Cao, Kang Liu, and Jun Zhao. 2025. Large language models for planning: A comprehensive and systematic survey. *arXiv preprint arXiv:2505.19683*.
- Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, and 1 others. 2025. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657*.
- Wendi Cui, Jiaxin Zhang, Zhuohang Li, Hao Sun, Damien Lopez, Kamalika Das, Bradley A Malin, and Sricharan Kumar. 2025. Automatic prompt optimization via heuristic search: A survey. *arXiv preprint arXiv:2502.18746*.
- Aryan Gulati, Brando Miranda, Eric Chen, Emily Xia, Kai Fronsdal, Bruno de Moraes Dumont, and Sanmi Koyejo. 2024. Putnam-axiom: A functional & static benchmark for measuring higher level mathematical reasoning in llms. In *Forty-second International Conference on Machine Learning*.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, and 1 others. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6.
- Wenyue Hua, Lizhou Fan, Lingyao Li, Kai Mei, Jianchao Ji, Yingqiang Ge, Libby Hemphill, and Yongfeng Zhang. 2023. War and peace (waragent): Large language model-based multi-agent simulation of world wars. *arXiv preprint arXiv:2311.17227*.
- Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, and 1 others. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274.

715	Harold W. Kuhn. 1955. The hungarian method for the assignment problem . <i>Naval Research Logistics Quarterly</i> , 2(1-2):83–97.	771
716		772
717		773
718	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.	774
719		775
720		776
721		777
722		778
723	Hang Li, Tianlong Xu, Ethan Chang, and Qingsong Wen. 2025a. Knowledge tagging with large language model based multi-agent system. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 28775–28782.	779
724		780
725		781
726		782
727		783
728	Wenwu Li, Xiangfeng Wang, Wenhao Li, and Bo Jin. 2025b. A survey of automatic prompt engineering: An optimization perspective. <i>arXiv preprint arXiv:2502.11560</i> .	784
729		785
730		786
731		787
732	Shuhang Lin, Wenyue Hua, Lingyao Li, Che-Jui Chang, Lizhou Fan, Jianchao Ji, Hang Hua, Mingyu Jin, Jiebo Luo, and Yongfeng Zhang. 2024. BattleAgent: Multi-modal dynamic emulation on historical battles to complement historical analysis . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 172–181, Miami, Florida, USA. Association for Computational Linguistics.	788
733		789
734		790
735		791
736		792
737		793
738		794
739		795
740		796
741	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. <i>arXiv preprint arXiv:2307.03172</i> .	797
742		798
743		799
744		800
745		801
746	Hao Ma, Tianyi Hu, Zhiqiang Pu, Liu Boyin, Xiaolin Ai, Yanyan Liang, and Min Chen. 2024. Coevolving with the other you: Fine-tuning llm with sequential cooperative multi-agent reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 37:15497–15525.	802
747		803
748		804
749		805
750		806
751		807
752	Melissa Z Pan, Mert Cemri, Lakshya A Agrawal, Shuyi Yang, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Kannan Ramchandran, Dan Klein, and 1 others. 2025. Why do multiagent systems fail? In <i>ICLR 2025 Workshop on Building Trust in Language Models and Applications</i> .	808
753		809
754		810
755		811
756		812
757		813
758	Laurent Perron and Vincent Furnon. 2022. Google OR-Tools. https://developers.google.com/optimization .	814
759		815
760		816
761	Jianing Qiu, Kyle Lam, Guohao Li, Amish Acharya, Tien Yin Wong, Ara Darzi, Wu Yuan, and Eric J Topol. 2024. Llm-based agentic systems in medicine and healthcare. <i>Nature Machine Intelligence</i> , 6(12):1418–1420.	817
762		818
763		819
764		820
765		821
766	Raphael Shu, Nilaksh Das, Michelle Yuan, Monica Sunkara, and Yi Zhang. 2024. Towards effective genai multi-agent collaboration: Design and evaluation for enterprise applications. <i>arXiv preprint arXiv:2412.05449</i> .	822
767		823
768		824
769		825
770		826
	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. <i>Nature medicine</i> , 29(8):1930–1940.	
	Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O’Sullivan, and Hoang D Nguyen. 2025. Multi-agent collaboration mechanisms: A survey of llms. <i>arXiv preprint arXiv:2501.06322</i> .	
	Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. 2024a. Mixture-of-agents enhances large language model capabilities. <i>arXiv preprint arXiv:2406.04692</i> .	
	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024b. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):186345.	
	Zhexuan Wang, Yutong Wang, Xuebo Liu, Liang Ding, Miao Zhang, Jie Liu, and Min Zhang. 2025. Agentdropout: Dynamic agent elimination for token-efficient and high-performance llm-based multi-agent collaboration. <i>arXiv preprint arXiv:2503.18891</i> .	
	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.	
	David H Wolpert. 2021. What is important about the no free lunch theorems? In <i>Black box optimization, machine learning, and no-free lunch theorems</i> , pages 373–388. Springer.	
	David H Wolpert and William G Macready. 1997. No free lunch theorems for optimization. <i>IEEE transactions on evolutionary computation</i> , 1(1):67–82.	
	Jiarong Wu, Songqiang Chen, Jialun Cao, Hau Ching Lo, and Shing-Chi Cheung. 2025. Isolating language-coding from problem-solving: Benchmarking llms with pseudoeval. <i>arXiv preprint arXiv:2502.19149</i> .	
	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, and 1 others. 2023a. Autogen: Enabling next-gen llm applications via multi-agent conversation. <i>arXiv preprint arXiv:2308.08155</i> .	
	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kam-badur, David Rosenberg, and Gideon Mann. 2023b. Bloomberggpt: A large language model for finance. <i>arXiv preprint arXiv:2303.17564</i> .	
	Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is inevitable: An innate limitation of large language models. <i>arXiv preprint arXiv:2401.11817</i> .	

827 Yuwei Yan, Qingbin Zeng, Zhiheng Zheng, Jingzhe
828 Yuan, Jie Feng, Jun Zhang, Fengli Xu, and Yong Li.
829 2024. Opencity: A scalable platform to simulate ur-
830 ban activities with massive llm agents. *arXiv preprint*
831 *arXiv:2410.21286*.

832 Ye Ye. 2025. Task memory engine (tme): Enhancing
833 state awareness for multi-step llm agent tasks. *arXiv*
834 *preprint arXiv:2504.08525*.

835 Huizi Yu, Jiayan Zhou, Lingyao Li, Shan Chen, Jack
836 Gallifant, Anye Shi, Xiang Li, Wenyue Hua, Mingyu
837 Jin, Guang Chen, and 1 others. 2024. Aipatient: Sim-
838 ulating patients with ehra and llm powered agentic
839 workflow. *arXiv preprint arXiv:2409.18924*.

840 Bowen Zhang and Pengcheng Luo. 2025. Or-llm-agent:
841 Automating modeling and solving of operations re-
842 search optimization problem with reasoning large
843 language model. *arXiv preprint arXiv:2503.10009*.

844 Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister,
845 Rui Zhang, and Sercan Arik. 2024. Chain of agents:
846 Large language models collaborating on long-context
847 tasks. *Advances in Neural Information Processing*
848 *Systems*, 37:132208–132237.

849 Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng
850 Yang, Shuyi Guo, Zhe Wang, Zhenhailong Wang,
851 Cheng Qian, Xiangru Tang, Heng Ji, and 1 others.
852 2025. Multiagentbench: Evaluating the collabora-
853 tion and competition of llm agents. *arXiv preprint*
854 *arXiv:2503.01935*.

A Use of LLMs

During the development of this paper, we use transformer based large language models in the following aspects:

- Reference discovery: use the deep research tools from major providers to explore relevant work and literature.
- Code assistance: use coding agents to assist developing the code base of the current work.
- Grammar check: use LLMs to detect grammar errors in the drafty version of the paper, for better displaying our results.

B Appendix: abstraction of our methodology

In this process, we provide formal definitions of the terminologies we used in Section 3.

Definition B.1. A well-formulated task is a tuple $T = (I, O, R)$, consisting of

- Input domain I : an unambiguous description of all admissible inputs.
- Output co-domain O : an unambiguous description of all admissible outputs.
- Requirement relation $R \subset I \times O$: a relation such that for each input $x \in I$ it defines explicitly the subset $R(x) \subset O$ as the set of outputs that are considered correct.

Definition B.2. A workflow blueprint $\mathcal{B} = (\mathcal{T}, \mathcal{P})$ consisting of

- A finite set of well-formulated tasks $\mathcal{T} = \{T_1, \dots, T_n\}$.
- An orchestration protocol \mathcal{P} that specifies: (i) The control-flow graph that determines when each T_i is invoked. (ii) The data-dependency edges that map outputs of some tasks to inputs of others. (iii) Any global invariants, error-handling rules, or communication channels required to realize the objective of \mathcal{B} .

Definition B.3. Given a workflow blueprint $\mathcal{B} = (\mathcal{T}, \mathcal{P})$, a **decomposition** selects a task $T \in \mathcal{T}$ and replace it with a sub-blueprint $\mathcal{B}_T = (\mathcal{T}_T, \mathcal{P}_T)$ such that (1). Each task $T' \in \mathcal{T}_T$ is strictly simpler than T . (2). The composite protocol \mathcal{P}' , obtained by embedding \mathcal{P}_T in place of T inside \mathcal{P} , preserves all external interface of T . The result of the decomposition is a new blueprint $\mathcal{S}' = (\mathcal{T}', \mathcal{P}')$, where $\mathcal{T}' = (\mathcal{T} \setminus \{T\}) \cup \mathcal{T}_T$.

Definition B.4. Let \mathcal{M} be a set of LLM models. A well-formulated task T is said to be \mathcal{M} -tractable if a model inside \mathcal{M} satisfies the requirement relation R_T with high, empirically verified accuracy, after optional augmentations (e.g., chain-of-thought prompting, tool calls, self-reflection loops, or fine-tuning).

Definition B.5. Given a set of LLM models \mathcal{M} and a blueprint \mathcal{B} , an \mathcal{M} -tractable hierarchy is a sequence of decompositions

$$\mathcal{B} \xrightarrow{D_1} \mathcal{B}_1 \xrightarrow{D_2} \dots \xrightarrow{D_n} \mathcal{B}_n$$

such that each task in the terminal blueprint \mathcal{B}_n is \mathcal{M} -tractable in the sense of Definition B.4.

Definition B.6. Given a set of LLM models \mathcal{M} and an \mathcal{M} -tractable blueprint $\mathcal{B} = (\mathcal{T}, \mathcal{P})$, a **system instantiation** is to instantiate \mathcal{B} into a MAS in the following way.

- Create one agent A_i per task $T_i \in \mathcal{T}$, bundling augmentations with the agent.
- We implement the orchestration protocol \mathcal{P} as message-passing or function calls among agents, preserving data-dependencies and control flow.

C Appendix: Weighted No Free Lunch Theorem—Why Does Agent Design Fail?

In this section, we present a weighted version of the No-Free-Lunch theorem. As the motivation, current approaches to MAS design can often result in overly general solutions that may exhibit suboptimal performance on specific and complex tasks. This sub-optimality arises partially from a lack of domain-specific inductive bias. To formalize this, we present a weighted variant of the No Free Lunch (NFL) theorem. The following demonstration, leveraging a weighted variant of the No Free Lunch theorem, quantitatively illustrates how inductive bias tailored to the target domain enhances performance. That is, We present a formal proof showing that, under a non-uniform prior concentrated on a problem-specific subset of functions, a specialized learning algorithm achieves strictly lower expected risk than a general-purpose algorithm.

Note the NLF theorem has been known to research community for more than two decades. Here what we present is a modification of the standard statement to better fit for our discussion on the MAS. As we didn't find in literature the precise version of the NFL theorem as we stated below, we also present a proof for self-containedness. We do not claim any originality of the theorem and the proof.

Theorem C.1 (Weighted NFL). *Let X be a finite input domain, Y a finite label set, and $\mathcal{F} = Y^X$ the set of all functions $f: X \rightarrow Y$. Consider*

- a general algorithm A_0 with constant expected loss ε_0 on every $f \in \mathcal{F}$,
- a specialized algorithm A' satisfying

$$L(h_{A'}, f) \leq \begin{cases} \varepsilon_1, & f \in \mathcal{F}', \\ \varepsilon_2, & f \notin \mathcal{F}', \end{cases}$$

where $\varepsilon_1 < \varepsilon_0 < \varepsilon_2$, and

- a prior P with $P(f \in \mathcal{F}') = p$ and $P(f \notin \mathcal{F}') = 1 - p$.

If

$$p > \frac{\varepsilon_0 - \varepsilon_2}{\varepsilon_1 - \varepsilon_2},$$

then the expected risk of A' is strictly lower than that of A_0 , i.e.

$$R(A') < R(A_0).$$

Proof. By definition,

$$R(A_0) = \mathbb{E}_{f \sim P}[L(h_{A_0}, f)] = \varepsilon_0$$

$$R(A') = \mathbb{E}_{f \sim P}[L(h_{A'}, f)] = p\varepsilon_1 + (1 - p)\varepsilon_2.$$

Hence

$$\begin{aligned} R(A') < R(A_0) &\iff p\varepsilon_1 + (1 - p)\varepsilon_2 < \varepsilon_0 \\ &\iff p(\varepsilon_1 - \varepsilon_2) > \varepsilon_0 - \varepsilon_2, \end{aligned}$$

which rearranges to

$$p > \frac{\varepsilon_0 - \varepsilon_2}{\varepsilon_1 - \varepsilon_2}.$$

This completes the proof. \square

D Appendix: KSP and TAP description

In this appendix, we provide details about the KSP and TAP, including their problem description and algorithm based on which we design our MAS.

D.1 KSP Problem Formulation

The usual input of KSP involves a set of N items, whose items are characterized by pairs (w_i, v_i) of weights w_i and values v_i , as well as a capacity value W . The goal of KSP is to find a subset of items such that the total weight does not exceed the given capacity while the total value is maximized. Mathematically, we record information of items by two vectors, both of dimension N : a weight vector $\vec{w} = (w_1, \dots, w_N)$ and a value vector $\vec{v} = (v_1, \dots, v_N)$. We also introduce the set of state-vectors $\{0, 1\}^N$, whose elements are vectors $\vec{x} = (x_1, \dots, x_N)$ where entries x_i takes values between 0 and 1, indicating whether an item is chosen in a subset or not:

$$x_i = \begin{cases} 1 & \text{item } i \text{ is chosen} \\ 0 & \text{item } i \text{ is excluded} \end{cases}$$

Thus state vectors controls which items is in the chosen subset, and the inner product of \vec{x} with \vec{w} and \vec{v} then compute the total weight and total value for the given subset, respectively.

Given a weight vectors \vec{w} , a value vector \vec{v} , and the capacity constraint W , the objective of the Knapsack problem then can be formulated as finding the following (optimal) value:

$$Z = \max\{ \vec{x} \cdot \vec{v} \mid \vec{x} \in \{0, 1\}^N, \vec{x} \cdot \vec{w} \leq W \}.$$

Here the maximal value is taken over all state vectors (or equivalently, all subsets of items) satisfying the constraint that the total weight $\vec{x} \cdot \vec{w}$ not exceeding the capacity W .

This version of the problem, where each item can either be fully included or not at all, is commonly known as the 0/1 KSP.

D.2 KSP Problem Solution

A classic approach to the Knapsack Problem iteratively enumerates all feasible states—a dynamic-programming strategy first introduced by Bellman (Bellman, 1957). A feasible state can be defined by a pair $(current_weight, current_value)$ representing the accumulated weight and value of a set of items selected so far, such that $current_weight \leq W$. We can describe the algorithm in the form of mathematical induction. We start with the initial set of feasible states $S_0 = \{(0, 0)\}$, representing an empty set of chosen items. We then add items in to form a set S_k from S_{k-1} inductively, with capacity being aware: for each k , assuming that S_{k-1} has been constructed, then we add the pair (w_k, v_k) to all items in S_{k-1} to form a new set S_{add} :

$$S_{add} = \{(w + w_k, v + v_k) \text{ for all } (w, v) \in S_{k-1}\}.$$

Then, we trim the set according to the capacity:

$$S_{trimmed} = \{(w, v) \in S_{add} \mid w \leq W\}.$$

Note this also removes all repetitive states in the set. Finally we take the union of the two intermediate sets to create S_k :

$$S_k = S_{k-1} \cup S_{trimmed}.$$

The inductive step terminate when we have run through all items and obtaining the final set S_N , we pick the element in S_N with maximal value, as the solution to the KSP. Explicitly,

$$Z = \max_{(w,v) \in S_N} v$$

D.3 TAP Problem Formulation

TAP seeks to optimally assign a set of N resources (agents or workers) to N tasks, where each potential assignment incurs a specific cost. With the constraint that each resource can only be assigned to one unique task, the objective of TAP is to find an assignment covering all tasks that minimizes the total cost. The resource-task specific cost is recorded in an $N \times N$ matrix C , where the entry C_{ij} represents the cost associated with assigning resource i to task j , for $i, j \in \{1, 2, \dots, N\}$.

To formally define the problem, we introduce a set \mathfrak{S}_N which can be described in either one of the following three equivalent ways:

- The group of automorphisms of the set $\underline{N} = \{1, 2, \dots, N\}$.
- The set (or group) of bijections from the set \underline{N} to itself.
- The set of all permutations involving N elements.

Note elements in \mathfrak{S}_N convey the idea that each resource is assigned to a unique task. Now, given the $N \times N$ cost matrix C , the objective of the TAP is to find the following (optimized) value

$$Z = \max_{\sigma \in \mathfrak{S}_N} \sum_{i=1}^n C_{i\sigma(i)}.$$

Note when we treat $\sigma \in \mathfrak{S}_N$ as a (bijective) map from $\underline{N} = \{1, 2, \dots, N\}$ to itself, the notation $C_{i\sigma(i)}$ represents the entry on the i -th row and $\sigma(i)$ -th column of the cost matrix C .

D.4 TAP Problem Solution

The typical solution for TAP is using the Hungarian algorithm (Kuhn, 1955), which provides a polynomial-time method to find the objective value Z . We summarize the algorithm as follows.

Step 1. Row Reduction. For each row, we find the minimal element in the row and subtract it from all entries in the row, creating at least one zero on each row. Mathematically, starting from the original cost matrix $C_{N \times N}$, we create a new reduced matrix C' such that for $i, j \in \{1, 2, \dots, N\}$, we have

$$C'_{ij} = C_{ij} - \min_{1 \leq k \leq N} C_{ik}$$

Step 2. Column Reduction. Similarly, we further reduce C' to C'' as follows: For each column, we find the minimal element in the column and subtract it from all entries in the column, guaranteeing at least one zero on each column. Mathematically, take

$$C''_{ij} = C'_{ij} - \min_{1 \leq k \leq N} C'_{kj}$$

Step 3. Find covering lines. We then find a smallest collection of rows and columns to cover all zeroes. Here smallest is the sense of the number of elements in the collection (of rows and column), over all possible such collections. If the size of this minimal collection, denoted by L , coincides with N , the size of the problem, then we skip Step 4 to enter the final stage of the algorithm.

Step 4. Matrix Improvement. However, if $L < N$, we need an improvement for the matrix C'' before looping back to Step 3: given the minimal collection of rows and columns from Step 3, we find the minimal value of all entries that are not covered, and then subtract this minimal value from all uncovered entries of C'' , and then add this minimal value to all entries of C'' that are covered twice, i.e., by both rows and columns. Let C''' be the resulting matrix.

Step 5. Assignment Identification. Once the condition $L = N$ is met, the final step is to identify the optimal assignment. This involves selecting a set of N independent zeros from the current matrix C''' , such that no two selected zeros share the same row or column. Each selected zero at position (i, j) corresponds to assigning agent i to task j . The total cost of this optimal assignment is then calculated by summing the costs from the original cost matrix C corresponding to selected zero positions.

A non-trivial fact guaranteed by the Hungarian algorithm is that, in Step 5 the collection of zeroes might not be unique, while different collections are deemed to result in the same total summation of corresponding entries in the original cost matrix C .

E Appendix: Academic Proofreading for Mathematical Papers

E.1 Model Performance

To further demonstrate the generality of KtR beyond structured optimization problems, we implement a LaTeX proofreading system for mathematical papers written by non-native authors. The proofreading task for mathematical papers requires: (i) exhaustively detecting grammar errors and providing detailed explanations, (ii) detecting notational inconsistencies, and (iii) being aware of mathematical terminologies, notations, and LaTeX commands. The same KtR controller framework now orchestrates a sentence-level “Grammar Agent” and an environment-level “Notation Agent,” demonstrating the framework’s adaptability to language-centric tasks.

Our implementation reveals two critical limitations in current LLMs. (i) Rule-following degradation: As input length increases, models become less likely to strictly adhere to system prompt instructions. (ii) Premature satisfaction: Models often terminate responses prematurely, even when token limits allow for further analysis. KtR addresses these limitations through principled task decomposition aligned with natural task boundaries: **(i) Grammar checking:** Grammatical correctness is determined at the sentence level. Thus, we decompose documents into individual sentences and employ a GPT-4o-mini agent for this language task. **(ii) Notation checking:** We determine that segmentation by math environments (definitions, statements, proofs) provides a balance between length and self-containedness. We employ an o3-mini agent for this reasoning task.

We run three experimental conditions: (i) zero-shot on the whole paper, (ii) zero-shot by sections, and (iii) our KtR-MAS approach. We configure parallel execution with 5 workers and tested on draft versions of research-level mathematics papers of various lengths (5, 18, 36, and 58 pages). Table 2 presents the grammar checking results across all methods and paper lengths, while Table 3 presents the notation checking results. Table 4 compares KtR-MAS against existing proofreading solutions.

Table 2: Grammar checking results across papers of varying lengths.

Paper	Pages	Method	Scope	Valid Suggestions	Time (s)	Cost (\$)
1	5	Zero-shot	Paper	12	23.35	0.00
		Zero-shot	Section	22	17.29	0.00
		KtR-MAS	Sentence	28	58.44	0.02
2	18	Zero-shot	Paper	15	30.66	0.00
		Zero-shot	Section	55	48.57	0.01
		KtR-MAS	Sentence	151	307.05	0.07
3	36	Zero-shot	Paper	19	43.69	0.01
		Zero-shot	Section	96	66.09	0.01
		KtR-MAS	Sentence	186	328.66	0.09
4	58	Zero-shot	Paper	20	37.60	0.01
		Zero-shot	Section	173	118.65	0.02
		KtR-MAS	Sentence	372	610.49	0.17

Table 3: Notation checking results across papers of varying lengths

Paper	Pages	Method	Scope	Critical	Advisory	Time (s)	Cost (\$)
1	5	Zero-shot	Paper	0	0	34.82	0.02
		Zero-shot	Section	1	4	22.04	0.03
		KtR-MAS	Environment	3	6	53.23	0.13
2	18	Zero-shot	Paper	3	0	26.50	0.04
		Zero-shot	Section	9	8	44.06	0.11
		KtR-MAS	Environment	19	69	261.68	0.70
3	36	Zero-shot	Paper	2	0	37.12	0.05
		Zero-shot	Section	12	8	89.48	0.15
		KtR-MAS	Environment	30	59	281.02	0.75
4	58	Zero-shot	Paper	0	2	36.14	0.07
		Zero-shot	Section	21	8	68.21	0.20
		KtR-MAS	Environment	63	85	439.19	1.22

Table 4: Comparison with existing proofreading methods.

Metric	Rule-based (GitHub)	AI-based (GitHub)	Writefull (Commercial)	KtR-MAS
Cost	Free	Token-based	Free/\$21 mo.	\$1.39
Grammar check	Yes	Limited ^a	Exhaustive	Exhaustive
Notation check	No	—	No	Yes
Explanations	No	—	Limited ^b	Yes
Math awareness	No	Varies	Rare ^b	Yes

^aProject designs resemble our zero-shot baselines. ^bFor free version; paid version not fully tested.

Our key findings are summarized below. **(i) Scalability:** KtR-MAS maintains effectiveness across document lengths, finding $18.5\times$ more valid grammar suggestions than zero-shot on the whole paper for the largest document, and $2\text{--}3\times$ more than zero-shot by sections. **(ii) Cost-effectiveness:** At \$1.39 for a 58-page paper with less than 20 minutes of wall-clock time, KtR-MAS is practical for real-world applications. This additional validation demonstrates that KtR’s principled decomposition approach generalizes beyond structured optimization problems to open-domain language tasks. The framework’s

ability to maintain performance while scaling to larger, more complex documents, combined with its cost-effectiveness and unique capabilities, supports our claim of broad applicability. The proofreading results show that KtR can handle ambiguous, context-dependent tasks requiring deep understanding of domain-specific conventions.

E.2 Cost and Latency Trade-offs

In addition, we incorporate the proofreading case study results to provide a concrete view of cost-versus-latency for our KtR-driven agent framework. All wall-clock times are measured end-to-end on a single GPU and therefore already include controller overhead, message routing, and the five concurrent worker agents; no post-hoc batching or log pruning is performed. Table 5 presents the cost and latency comparison between zero-shot (by sections) and KtR-MAS approaches for grammar checking, while Table 6 presents the corresponding analysis for notation checking.

Table 5: Cost and latency comparison for grammar checking.

Paper	Pages	Method	Total Cost (USD)	Time Latency (s)
1	5	Zero-shot	0.00	23.35
		KtR-MAS	0.02	58.44
2	18	Zero-shot	0.00	30.66
		KtR-MAS	0.07	307.05
3	36	Zero-shot	0.01	43.69
		KtR-MAS	0.09	328.66
4	58	Zero-shot	0.01	37.60
		KtR-MAS	0.17	610.49

Table 6: Cost and latency comparison for notation checking.

Paper	Pages	Method	Total Cost (USD)	Time Latency (s)
1	5	Zero-shot	0.02	34.82
		KtR-MAS	0.13	53.23
2	18	Zero-shot	0.11	44.06
		KtR-MAS	0.70	261.68
3	36	Zero-shot	0.15	89.48
		KtR-MAS	0.75	281.02
4	58	Zero-shot	0.20	68.21
		KtR-MAS	1.22	439.19

For the energy/CO₂ budget, using the commonly cited estimate of 0.3 Wh per 1,000 tokens for GPT-class inference, the 328,933 output tokens in the 58-page case translate to approximately 0.10 kWh, again underscoring that optimization happens through low-power, small-model calls. The zero-shot results presented above use section-level decomposition as the baseline, representing a stronger comparison than whole-paper zero-shot inference. We also observe that the zero-shot time latency for the 58-page paper appears anomalously low. We note that this is indeed what is recorded by the test program. Since we segment papers into sections and configure parallel execution with 5 workers, time latency may be affected by paper structure, API connection speeds, and response variability.

F Ground-truth Data Preparation

We utilize Google OR-Tools (Perron and Furnon, 2022) to generate optimal solutions—serving as ground-truth datasets—for both problem scenarios. OR-Tools is a widely adopted open-source software suite developed by Google for solving combinatorial optimization problems. It is renowned for its efficiency and reliability in addressing NP-hard challenges through advanced optimization algorithms. The suite is distributed under the permissive Apache License 2.0, allowing unrestricted use, modification, and distribution (Perron and Furnon, 2022).

For KSP, we generate random instances by assigning weights and values to items along with a maximum capacity constraint. Optimal solutions are then computed using OR-Tools' dynamic programming approach. For TAP, we similarly generate random cost matrices that represent the cost of assigning workers to tasks. Optimal assignments are obtained using the Hungarian algorithm as implemented in OR-Tools, which efficiently minimizes the total assignment cost.

G Appendix: Prompt gallery

Note that all prompts we presented in the following, except for the self-check prompt for Trimmer Agent in KSP problem, are the system prompt for agents. The user prompt will only contain the precise problem to be handled by the agent in the form specified by the prompt.

G.1 KSP prompts

G.1.1 Prompt for zero shot

```
You are an expert in the field of Knapsack Problem.

You are given a Knapsack Problem in the json format, of the following form:
{
  "id" : str,
  "items" : list of pairs of integers,
  "capacity" : int
}

Each pair in the list is a pair of integers of the form [weight, value], i.e., the
first entry is the weight and the second entry is the value.

Your task is to solve the Knapsack Problem and provide the optimal solution. That
is, you need to find a subset of the pairs that maximizes the total value,
subject to the constraint that the total weight of the subset is less than or
equal to the capacity.

Please think step by step when solving the problem.

You need to return the optimal solution in the following json format:
{
  "max_value" : int,
}

Please only return the json format, nothing else.
```

G.1.2 Prompt for Worker Agent

```
You are a key member of a multi-agent team collaboratively solving the Knapsack
Problem. Your specific role is the Worker, responsible for performing
mathematical computations for the team.

You will receive input in the following JSON format:
{"c_list": [[int, int], ...], "s_item": [int, int]}
Each pair within 'c_list' contains two integers.

Your task is to:
- Add 's_item' to each pair in 'c_list' entry-wise. For instance, if a pair in
  'c_list' is '[2, 5]' and 's_item' is '[3, 4]', the result should be '[2+3, 5+4]
  = [5, 9]'.

To ensure accuracy:
- Proceed systematically, applying step-by-step reasoning.
- Carefully perform each addition individually for all pairs provided in the list.

Your response must strictly follow this JSON format:
{"n_list": [[int, int], ...]}

Return only the specified JSON object without any additional commentary or text.
```

G.1.3 Prompt for Trimmer Agent

You are a key member of a multi-agent team collaboratively solving the Knapsack Problem. Your specific role is the Trimmer, responsible for trimming the list based on the given capacity constraint.

You will receive input in the following JSON format:
{"n_list": [[int, int], ...], "capacity": int}

Each pair within 'n_list' contains two integers: the first integer represents the weight, and the second integer represents the value.

Your task is to:

- Carefully analyze each pair in the provided list.
- Remove all pairs whose weight (the first integer) strictly exceeds the specified capacity.
- If identical pairs appear multiple times, retain only one instance of each.

To ensure accuracy:

- Proceed systematically, applying step-by-step reasoning.
- Verify each pair carefully against the capacity constraint.

Your response must strictly follow this JSON format:
{"t_list": [[int, int], ...]}

Return only the specified JSON object without any additional commentary or text.

G.1.4 Prompt for Reporter Agent

You are a key member of a multi-agent team collaboratively solving the Knapsack Problem. Your specific role is the Reporter, responsible for determining and clearly reporting the final answer based on the provided information.

You will receive input in the following JSON format:
{"c_list": [[int, int], ...]}

Each pair within 'c_list' contains two integers: the first integer represents the weight, and the second integer represents the value.

Your task is to carefully analyze this list, identify the pair with the maximal value (the second integer in each pair), and report only that maximal value. If the list is empty, then report the maximal value as 0.

To ensure accuracy:

- Proceed systematically, applying step-by-step reasoning.
- Carefully examine every pair in the provided list.

Your response must strictly follow this JSON format:
{"max_value": int}

Return only the JSON object as specified above, without any additional commentary or text.

G.1.5 Self-check prompt for Trimmer Agent

To better fulfill your task, please conduct a double check on the result you just provided. If your answer is already correct, please confirm by copying the last output.

When double check, please pay attention to the following typical types of mistakes:

In particular, please check if you made any typical mistakes as listed below:

1. If you added in a pair that is not in the original n_list.
2. If there is still a pair in the t_list that still exceeds the capacity.
3. If there is a pair in n_list that does not exceed the capacity but is not in the t_list.

If you found any errors, please create a corrected answer.

In either case, please follow the format requirement of the output.

G.2 TAP prompts

G.2.1 Prompt for zero shot

You are an expert in solving the Assignment Problem. In the assignment problem, there are n workers and n jobs. Each worker has a cost of assigning to each job. Each worker can only be assigned to one job. Your task is to find the optimal assignment of workers to jobs that minimizes the total cost.

You are given the problem in the following json format:

```
{
  "id" : str,
  "cost_matrix" : list of lists of integers
}
```

The cost matrix is a square matrix of size $n \times n$, where n is the number of workers and jobs, in the form of a nested list `[[int, int, ...], [int, int, ...], ...]`. The (i, j) th entry of the matrix represents the cost of assigning the i th worker to the j th job.

Your task is to find the optimal assignment of workers to jobs that minimizes the total cost.

Please think step by step when solving the problem.

You need to return the optimal assignment in the following json format:

```
{
  "optimal_cost" : int
}
```

Please only return the json format, nothing else.

G.2.2 Prompt for Row Reducer Agent

You are given a matrix in the following json format:

```
{
  "matrix" : list of lists of integers
}
```

The matrix is in the form of a nested list `[[int, int, ...], [int, int, ...], ...]`.

Your task is to reduce the matrix by subtracting the minimum value of each row from all the elements in that row.

Please think step by step when solving the problem:

Step 0: Work on one row at a time.

Step 1: Find the minimum value of the row.

Step 2: Subtract the minimum value of the row from all the elements in that row.

Step 3: Return the reduced matrix in the following json format:

```
{"reduced_matrix" : list of lists of integers}
```

Please only return the json format, nothing else.

G.2.3 Prompt for Column Reducer Agent

You are given a matrix in the following json format:

```
{
  "matrix" : list of lists of integers
}
```

The matrix is in the form of a nested list `[[int, int, ...], [int, int, ...], ...]`.

Your task is to reduce the matrix by subtracting the minimum value of each column from all the elements in that column.

Please think step by step when solving the problem: 1268
 Step 0: Work on one column at a time. 1269
 Step 1: Find the minimum value of the column. 1270
 Step 2: Subtract the minimum value of the column from all the elements in that 1271
 column. 1272
 Step 3: Return the reduced matrix in the following json format: 1273
 1274
 {"reduced_matrix" : list of lists of integers} 1275
 1276
 Please only return the json format, nothing else. 1277
 1278

G.2.4 Prompt for Zero Seeker Agent 1280

You are given a problem in the following json format: 1281
 1282
 { 1283
 "matrix" : list of lists of integers 1284
 } 1285
 1286
 The matrix is in the form of a nested list [[int, int, ...], [int, int, ...], ...]. 1287
 1288
 Your task is to find a smallest collection of rows and columns of the matrix, such 1289
 that any zeroes in the matrix is contained in a chosen row or column. Small 1290
 means the sum of the sizes of the row and column collections is the smallest 1291
 possible. 1292
 1293
 Please think step by step when solving the problem, and return your response in the 1294
 following json format: 1295
 1296
 {"collum_collection" : [int, int, ...], "row_collection" : [int, int, ...]} 1297
 1298
 The integers in the collum_collection and row_collection are the indices of the 1299
 rows and columns that you choose. 1300
 1301
 Please only return the json format, nothing else. 1302
 1303

G.2.5 Prompt for Matcher Agent 1305

You are given a matrix in the following json format: 1306
 1307
 { 1308
 "matrix" : list of lists of integers 1309
 } 1310
 1311
 The matrix is in the form of a nested list [[int, int, ...], [int, int, ...], ...]. 1312
 1313
 Your task is to find the largest collection of zeroes in the matrix, such that no 1314
 two zeroes are in the same row or column. 1315
 1316
 Please think step by step when solving the problem, and return your response in the 1317
 following json format: 1318
 1319
 {"largest_collection" : [[int, int], [int, int], ...]} 1320
 1321
 The list of pairs of integers is in the form of [[row_index, column_index], 1322
 [row_index, column_index], ...]. 1323
 1324
 Please only return the json format, nothing else. 1325
 1326

G.2.6 Prompt for Painter Agent 1328

You are given a problem in the following json format: 1329
 1330
 { 1331
 "matrix" : list of lists of integers 1332
 "collection" : list of lists of integers 1333
 } 1334
 1335

1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363

The matrix is in the form of a nested list `[[int, int, ...], [int, int, ...], ...]`.

The collection is in the form of a nested list `[[int, int], [int, int], ...]`.

Your task is to find a smallest collection of rows and columns of the matrix, such that any zeroes in the matrix is contained in a chosen row or column. Small means the sum of the sizes of the row and column collections is the smallest possible.

To assist you, you are provided with a collection of zeroes in the input json format. The collection contains the positions of a maximal collection of zeroes in the matrix, such that no two zeroes are in the same row or column.

Please use this collection of zeroes to find the rows and columns as desired. More precisely, you should first choose one row or column for each zero in the collection, such that the chosen rows and columns cover as much of the zeroes in the matrix as possible. Then add in more rows or columns if needed.

Please think step by step when solving the problem, and return your response in the following json format:

```
{"collum_collection" : [int, int, ...], "row_collection" : [int, int, ...]}
```

The integers in the `collum_collection` and `row_collection` are the indices of the rows and columns that you choose.

Please only return the json format, nothing else.

1365

G.2.7 Prompt for Normalizer Agent

1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393

You are given a problem in the following json format:

```
{  
  "matrix" : list of lists of integers  
  "collumn_collection" : list of integers  
  "row_collection" : list of integers  
}
```

The matrix is in the form of a nested list `[[int, int, ...], [int, int, ...], ...]`. The `collumn_collection` and `row_collection` are the indices of some selected rows and columns that covers all the zeroes in the matrix.

Your task is the following:

1. Find the minimal value in the matrix that is not covered by the selected rows and columns.
2. If this value is 0, return the original matrix.
3. If this value is not 0, subtract this value from all uncovered entries in the matrix.
4. For the entries that covered by both a selected row and a selected column, add this value to the entries.
5. For the entries that are covered by a selected row or column, but not both, do nothing.
6. Please return the updated matrix in the following json format:

```
{"normalized_matrix" : list of lists of integers}
```

Please only return the json format, nothing else.

1395

G.2.8 Prompt for Reporter Agent

1396
1397
1398
1399
1400
1401
1402
1403
1404

You are given a problem in the following json format:

```
{  
  "matrix" : list of lists of integers  
  "collection" : list of lists of integers  
}
```

The matrix is in the form of a nested list `[[int, int, ...], [int, int, ...], ...]`.

The collection contains a set of entries of the matrix in the form of `[[row_index, column_index], [row_index, column_index], ...]`.

Your task is the following:

1. Sum up the values of all the entries in the collection.
2. Return the total value in the following json format:

```
{"total_value" : int}
```

Please only return the json format, nothing else.

1405
1406
1407
1408
1409
1410
1411
1412
1413
1414