**⑥ ChatGPT**

# Optimization-Based Coordination of Multi-Agent AI Systems: Comprehensive Survey and Outlook

## Part 1: Technical Survey

### 1. Agent-Task Assignment Optimization

**Quadratic Assignment Problem (QAP) – Classical Formulation:** The QAP models the assignment of $n$ agents to $n$ tasks with pairwise interaction costs, seeking a permutation that minimizes total cost. It is NP-hard, and exact solvers (branch-and-bound, cutting planes, linearizations) can only handle small sizes (up to $n\approx 30$). Modern MILP solvers (Gurobi, CPLEX) can encode QAP but struggle as $n$ grows. Specialized branch-and-bound algorithms (e.g. Hahn & Grant, 2001) and cutting-plane relaxations (e.g. based on semidefinite programming) push exact methods slightly further, but remain exponential in worst-case. **State-of-the-art exact QAP solvers in 2024** include parallel branch-and-bound implementations and quantum-inspired optimizers (e.g. Toshiba's SQBM+ or LightSolver) which reported solving instances up to $n\approx 60$ within 60 seconds. However, most practical instances still require heuristic or metaheuristic approaches.

**Heuristics & Metaheuristics for QAP:** Due to the size limitations of exact solvers, numerous heuristic algorithms have been developed for QAP. *Local search* (swap-based neighborhood) was explored from the 1960s, with modern variants like **breakout local search** improving efficacy. Among *trajectory methods*, **tabu search** remains particularly effective – Taillard's 1991 robust tabu search is *still* considered one of the most successful QAP heuristics in terms of simplicity and performance [1] . Other single-solution methods include simulated annealing and iterated local search. *Population-based metaheuristics* have achieved the best known results on large QAPs: **genetic algorithms (GAs)** and hybrid GA variants (memetic algorithms) are consistently top-performers. Modern hybrids combine GAs with tabu or other local searches for intensification. For example, a 2024 hybrid GA + hierarchical iterated tabu search (Misevičius *et al.* 2024) achieved *three new best-known solutions* on challenging QAP instances [2] , demonstrating that careful hybridization can advance the state of the art. Many nature-inspired algorithms (ant colony, particle swarm, etc.) have also been adapted to QAP, though they often require hybridization with problem-specific local searches to be competitive. **Best solvers in practice (2024)** often use hybrid tabu/GAs or advanced tabu alone. Notably, Misevičius's 2024 algorithm found >90% of runs yielding optimal or near-optimal solutions on instances up to size 300+, indicating the power of modern hybrids [2] . Open-source implementations for QAP heuristics exist (e.g. **Matlab/Ejor** repository for tabu search, or Python code from recent papers), but many researchers still implement custom solvers due to problem-specific tweaks. Some learning-based approaches are emerging (see below), but classical metaheuristics remain dominant for QAP.

**Dynamic and Contextual Assignment Extensions:** In an AI research system, the "cost" of assigning a particular agent to a task might be context-dependent (agent's current load, task complexity, synergy with other agents) – a scenario not captured by a static QAP cost matrix. *Online QAP* (dynamic assignments over time) is an open challenge. One approach is to treat it as a sequence of QAP instances and use **warm-start heuristics**: e.g. start each new assignment from the previous solution and apply local re-optimization (tabu

search can be re-seeded with prior solutions). Some literature on *dynamic facility layout* (the QAP variant for plant layout changing over time) suggests heuristic strategies for slowly evolving cost matrices, but comprehensive online QAP solvers are lacking. Recent work in multi-agent systems emphasizes the need to handle **non-stationary agent performance** and drift. For instance, Tarasova *et al.* (2025) present a decentralized task allocation where agents continually update their task utility estimates via stochastic approximation and consensus. Their approach effectively *learns* assignment decisions online, addressing issues like agent fatigue or skill drift. This hints that **learning-based cost models** could be useful: e.g. using reinforcement learning to predict an agent-task compatibility score given system state. Early steps in this direction include training ML models to predict assignment quality; one could imagine a Graph Neural Network that inputs the system state and outputs a "cost matrix" for assignment, updated each iteration. However, integrating learned models into combinatorial optimization is non-trivial – ensuring the learned cost function yields good global assignments is an open research gap (see Section 7 Gap 1).

**Learning-Driven Assignment Methods:** A few recent studies try to marry learning with assignment solving. *Learning to solve QAP:* Liu *et al.* (2022) introduced an RL method (RGM) that learns construction heuristics for QAP, successfully solving all QAPLIB instances without using exact solvers [3]. More recently, Transformers have been trained to iteratively improve QAP solutions. For example, a 2024 approach uses a transformer policy (SAWT) that learns to perform local swap improvements, guided by a reward signal tied to assignment cost. These learning-based solvers can generalize to new instances and sometimes approach heuristic solver quality, though they may not consistently hit optimal solutions. Another angle is **contextual bandits for assignment** – an agent could treat each assignment decision as an arm selection, with context being task features. However, this is complex due to the combinatorial coupling of decisions (this essentially reduces to learning a good assignment function, which is as hard as the original problem). Overall, learning-based assignment is an emerging area with a handful of proof-of-concept results, but no off-the-shelf implementations yet; it remains a promising direction for adapting to dynamic contexts.

**Auction & Market-Based Task Allocation:** In multi-robot and multi-agent communities, *auction algorithms* are a popular decentralized way to assign tasks. Each agent "bids" for tasks based on its own utility, and a coordinator (or distributed consensus) assigns tasks to maximize total utility. **Single-item auctions** (each task auctioned independently) are simple and used in systems like the Contract Net Protocol (CNP). They are efficient but may not capture complex task couplings. **Combinatorial auctions** (agents bid on bundles of tasks) can handle synergies but are NP-hard to solve optimally; approximate winner-determination algorithms or iterative auctions are used in practice. Auction-based allocation has been successfully applied in multi-UAV task assignment, satellite scheduling, etc., because it naturally accommodates self-interested agents and partial observability. For our cooperative AI agents, auctions could be an *optimization heuristic*: e.g. each agent has an internal score for handling a task given the current state, and tasks are assigned to the highest bidder. This can approximate a greedy assignment and can be fast and adaptive. However, **auction methods assume static agent capabilities** during an auction round and often rely on broadcasts [4]. In dynamic environments, their performance can degrade if agent models or bids don't update quickly when conditions change [4]. There is research on *dynamic auctions* where tasks arrive online or agent bids are updated periodically, as well as *market-based mechanisms* where a virtual currency or prices guide the allocation. These approaches ensure a form of incentive compatibility (if agents were truly self-interested), but in our system all agents ultimately share the same goal, so incentive alignment is less of a concern. **Practical note:** Auction algorithms are easy to implement (especially single-task auctions) and could serve as a baseline for task assignment. Open-source: the **ROSWild** multi-robot framework and MATLAB's auction solvers (e.g. Murphey's auction code) can be adapted. Tools like Google OR-Tools can also solve assignment

via the Hungarian algorithm (for linear assignment) or even solve small combinatorial auction allocations via MIP.

**Multi-Objective and Constrained Assignment:** Often we have multiple criteria for "optimal" assignments – e.g. maximize research quality, minimize total runtime, and ensure critical roles (validators) are assigned. A *multi-objective assignment problem* can be formulated, where each assignment has several cost metrics. Solutions can be evaluated in a Pareto-optimal sense. Common techniques include **scalarization** (weighting objectives into one cost) and **$\epsilon$-constraint** methods (optimize one objective while enforcing thresholds on others). For example, one could minimize execution time subject to achieving at least a certain quality score. Researchers have developed methods to enumerate Pareto-optimal assignments for small problems, but this becomes intractable as size grows. Metaheuristics like NSGA-II (a genetic algorithm for multi-objective optimization) can be applied to assignments by encoding permutations and using crossover/mutation, yielding an approximation of the Pareto front. Some works also use **goal programming** to handle soft constraints (e.g. penalty for violating "Validator must be included"). **Constraint programming** is another approach: one can encode "agent X must get a task" or "at least one validator in assignment" as constraints and let a CP-SAT solver find any feasible assignment that also optimizes a main objective. OR-Tools CP-SAT, for instance, could handle this for moderate sizes (tens of agents). In summary, multi-objective assignment is typically handled by either scalarizing (with careful tuning of weights) or using evolutionary multi-objective algorithms to sample trade-off solutions. There is no one-size-fits-all solver – the approach depends on the priorities (e.g., if one objective like quality is paramount, treat others as constraints). Identifying **Pareto-optimal agent-task matchings** could be a valuable analysis: e.g. one assignment might be fastest, another highest-quality; understanding that frontier helps decide a suitable compromise. Currently, this is mostly done in an ad-hoc way; integrating multi-objective optimization deeply into agent coordination (perhaps via a meta-controller that adjusts weights based on user preferences) is a potential research direction (Gap 1 and Gap 3 in Section 7).

## 2. Workflow Routing & DAG Scheduling

**DAG Scheduling Fundamentals:** Once tasks (or sub-problems) are defined and assigned to agents, we must schedule their execution respecting dependencies. This often forms a Directed Acyclic Graph (DAG) of tasks, where nodes are subtasks and edges represent precedence (e.g. "Critic must follow Designer"). *Classic DAG scheduling* (from parallel computing and project management literature) focuses on minimizing total makespan (completion time) or meeting deadlines. Techniques like the **Critical Path Method (CPM)** find the longest dependency path to identify bottlenecks, and heuristics then prioritize scheduling of tasks on that critical path to reduce overall time. For example, if Designer → Critic → Validator is the critical chain, we ensure minimal waiting between them. **List scheduling** is a greedy approach: assign each task a priority (e.g. by earliest start or longest processing time) and then repeatedly schedule the ready task with highest priority on the first available resource. Variants like **HEFT (Heterogeneous Earliest Finish Time)** are commonly used for DAG scheduling on heterogeneous machines – HEFT computes a rank for each task (based on downstream path lengths) and then schedules tasks in descending rank order on the resource that yields earliest finish time. These heuristics have been shown effective for static DAGs. In our context, tasks are handled by specialized agents (which can be seen as resources with different speeds/skills). A similar approach can be taken: e.g. prioritize critical subtasks (like those gating the research outcome) and assign them to the fastest capable agent first. If certain agents can perform tasks in parallel, scheduling must account for resource contention (e.g. two heavy tasks shouldn't go to the same agent simultaneously if it's single-threaded). In summary, borrowing from PERT/CPM, we'd first map out the **workflow graph**, identify parallelism vs sequential parts, and then apply a heuristic to order execution.

**Sequential vs. Parallel Execution Trade-offs:** Unlike a static pipeline, an AI research workflow might allow some agents to work concurrently on independent subtasks (e.g. multiple "Literature Scout" agents searching papers in parallel). Deciding which agents run in parallel vs in sequence is itself an optimization. Running everything in parallel maximizes speed but could overload resources or lead to incoherent results (like too many divergent ideas). Running strictly sequentially (Designer then Critic then Refactorer, etc.) ensures a disciplined hand-off but may be slow and under-utilize available computation. This becomes a **makespan vs. quality** trade-off (since parallel exploration can improve quality up to a point by exploring alternatives, but too much parallelism might be wasteful). The scheduling problem is thus multi-objective: minimize time, but also ensure sufficient thoroughness. A possible strategy: allocate a *parallel pool* of agents for exploratory tasks (like multiple hypothesis generation in parallel), then join results and proceed sequentially for critical evaluation tasks. There is a parallel to **pipeline parallelism vs. data parallelism** in ML – sometimes it's worth having redundancy (parallel agents working on same step to get diverse outputs). Existing scheduling algorithms can incorporate simple cost models (e.g. assign a higher "cost" to using an agent in parallel if it's already loaded) to gently encourage load balancing.

**Adaptive & Confidence-Aware Routing:** A unique aspect of an LLM-based multi-agent workflow is that *the path can depend on content quality*. For example, if the Designer's output is high confidence, maybe we can skip an extensive critique phase and go straight to validation, whereas if the output seems flawed, we route it to a Refactorer or request more scrutiny. This resembles a **conditional DAG or a flowchart** where node execution is contingent on runtime conditions. Traditional DAG scheduling assumes a fixed structure, but here we want dynamic decision points. One way to handle this is with *conditional branching* in workflows: e.g. use a "confidence score" from the Critic agent as a gating condition – if score > 0.9, skip directly to final validation; if < 0.5, perhaps loop back to the Designer for a second attempt; if in between, go to a Refactorer for improvement. Many **workflow orchestration tools do support conditional tasks** (for instance, Apache Airflow's `BranchPythonOperator` or Prefect's conditional logic allow skipping certain branches based on runtime criteria). Prefect 3.0 explicitly embraces *dynamic workflows* where you can include normal Python `if/else` logic to decide which tasks to run at runtime. This is crucial for confidence-aware routing. We could implement a controller agent that reads intermediate results and triggers the next appropriate agent(s) accordingly. The scheduling challenge becomes **stochastic scheduling** – the workflow is no longer deterministic, so we might model each branch with a probability (e.g. there's a 20% chance the Critic finds the design acceptable and we skip straight to final step). This becomes akin to *scenario-based scheduling*: plan for different possible execution paths. Algorithms from *stochastic project scheduling* (which often use Monte Carlo simulation or robust scheduling techniques) might be applicable. State-of-the-art research in *stochastic DAG scheduling* considers uncertain task durations rather than conditional tasks, but some techniques overlap – e.g. **dynamic scheduling policies** that make on-the-fly decisions based on observed task outcomes. In reinforcement learning terms, one could treat the workflow routing as an RL policy: an agent observes the current state (e.g. confidence scores, partial results) and decides which agent to invoke next. Indeed, there are early works like **Captain Agent (Song et al., 2024)** that attempt dynamic re-routing of LLM agents in response to partial results. Another example is **DyLAN (Liu et al., 2023)** which explored dynamic role assignment in multi-agent dialogue. However, these are relatively nascent frameworks and not general scheduling algorithms.

From a practical tool perspective, **Airflow, Prefect, and Kubeflow Pipelines** are widely used for orchestrating DAGs of tasks: - *Apache Airflow:* Great for reliable execution of static DAGs with scheduling, retries, etc. It has some support for branching (via conditional operators) but the DAG structure must be defined in advance. It doesn't natively support dynamically adding new tasks or changing the DAG after execution starts. You can implement a feedback loop in Airflow, but it's often clunky (e.g. task X triggers a

new DAG run). - *Prefect:* More dynamic and Pythonic. Prefect flows can include normal control flow; as noted, Prefect 2.0+ removed the strict DAG requirement, allowing dynamic branching and looping within a flow. Prefect also makes it easier to run flows locally or in various environments. For adaptive workflows, Prefect would be a better fit than Airflow because you can literally write `if score > 0.9: skip_task()` logic. In fact, Prefect's docs highlight skipping branches based on runtime info as a key feature. - *Kubeflow Pipelines:* Designed for ML workflows on Kubernetes. It's more static (like Airflow) but integrates with TensorFlow and Kubernetes jobs. It supports conditionals and loops to some extent (with static unrolling or via custom components). However, it's less commonly used for LLM agent orchestration – it's more for ML model training pipelines. - *Others:* There are research orchestration frameworks like **LangChain's LangGraph** which treat agent workflows as graphs with possible cycles and branching. These are discussed in Part 2.

**Quality-of-Service and Real-Time Considerations:** If we consider the research system responding to queries or deadlines, scheduling might need to ensure certain tasks complete by a time. Real-time scheduling algorithms (rate-monotonic, earliest-deadline-first) exist but assume periodic tasks, which might not directly apply here. However, if the system is interactive (say user waiting for a research summary), a soft deadline can be imposed. A possible strategy is **progressive refinement**: allocate a certain amount of time to the workflow and schedule agents such that a reasonable result is ready by the deadline, with any remaining time used for incremental improvements. This might mean scheduling a quick initial pass (Designer → Validator for a rough answer) in parallel with deeper analysis tasks that might finish later (Designer → Critic → Refine, etc.), so that at least some result is ready in time. This is analogous to anytime algorithms where you have a coarse result quickly and refine over time.

**Routing and Feedback Loops:** Workflows like Designer → Critic → Refactorer → (back to Critic...) inherently contain loops. This is *not* a DAG (it's a cyclic graph). Traditional scheduling doesn't cover loops since it assumes DAGs. But many multi-agent workflows will be iterative: e.g. generate hypothesis, critique, improve, loop until satisfied or out of time. Handling this requires either unrolling the loop a fixed number of times in a DAG (if an upper bound on iterations is known), or using a while-loop structure in a dynamic orchestrator (Prefect supports normal loops). Monitoring such loops requires **termination criteria** (like if improvement < ε or max 3 iterations). Reinsertion of tasks (like re-invoking Designer) can be done programmatically in frameworks like LangChain's **Agent loops** or by a master agent deciding to iterate. This is less about scheduling optimality and more about system control logic, but the *optimization question* is: how many iterations to allow, and how to detect diminishing returns automatically? That's an open area: meta-level policies could learn to stop early if further loops likely won't help (perhaps via a learned value of perfection vs. time trade-off).

**State of the Art – Stochastic/Robust Scheduling:** In cloud computing and workflow literature (2020-2024), there's interest in **robust scheduling** (accounting for variability) and even *self-healing workflows*. For example, schedulers that can re-order tasks on the fly when some fail or when resources become unavailable. Our scenario might see an agent "fail" (e.g. produce nonsense). A robust approach would detect this (validator flags it) and then perhaps route the task to a different agent or re-run. This is analogous to task retry policies in Airflow (which simply re-run the same task) but extended to *alternate agents* or methods. No off-the-shelf framework does this out-of-the-box; it would require custom logic (like "if Agent_A fails, use Agent_B and then continue").

**Research outlook:** The key challenge for scheduling in our system is **adaptivity**. Existing DAG scheduling algorithms give a baseline for static workflows. The novel part is incorporating quality signals and

branching. Likely we will design a custom orchestrator or use a hybrid of a workflow engine (for reliability) and an RL-based decision maker for routing. State-of-the-art pipeline tools can serve as a foundation (Prefect is promising), but we might extend them with a custom scheduler that queries agent confidence scores and decides which branch to activate next. In Section 7 (Gap 2) we discuss potential research into *confidence-aware DAG scheduling*, as currently few if any frameworks or algorithms explicitly handle that.

## 3. Resource Allocation & Meta-Learning

**Multi-Armed Bandit (MAB) Algorithms for Agent Resource Allocation:** When multiple agents (or strategies) compete for limited resources (CPU time, API calls), we can cast it as a multi-armed bandit problem: each agent is an "arm" that we pull by allocating resources, and we get a reward (e.g. task success or quality). The goal is to allocate resources to maximize cumulative reward (overall system performance) under constraints. Classic bandit algorithms provide principled exploration vs. exploitation: - **Epsilon-Greedy:** Simple heuristic where most of the time you allocate to the best-known agent, but a fraction $\varepsilon$ of the time you try others. Easy to implement but not optimal. - **Upper Confidence Bound (UCB):** A family of algorithms (Auer *et al.*, 2002) that compute an optimism-adjusted score for each arm: roughly, estimated mean performance + uncertainty margin. UCB1, for example, picks the agent with the highest $\text{mean}_i + c\sqrt{\frac{\ln t}{n_i}}$ at round $t$, balancing exploration (favoring arms with few pulls $n_i$) and exploitation. UCB has strong regret guarantees and is simple. There are variants like UCB-Tuned (which adapts the confidence interval to variance) and Bayesian UCB. - **Thompson Sampling (TS):** A Bayesian approach where we maintain a probability distribution (posterior) for each agent's performance and sample from it to decide allocation. For example, if we model an agent's success probability with a Beta distribution, we sample a success probability from each agent's Beta posterior and allocate to the one with highest sample. Thompson Sampling tends to perform very well in practice and is easy to implement. It can also handle cases where rewards are not simply success/failure by using appropriate distributions (Gaussian TS for continuous rewards, etc.). Libraries like `mabwell` or `MABWiser` (Python) implement various bandits including TS and UCB. - **Contextual Bandits:** If each decision has a context (features), e.g. the nature of the task or current system state, contextual bandits use those features to select arms. Algorithms like LinUCB (which learns a linear model of reward given context) or neural bandits (which use a neural net to predict rewards) could allow the system to allocate resources adaptively: e.g. "for coding tasks, agent A tends to do better, whereas for writing tasks, agent B is better" – the context (task type) informs the decision. This fits our scenario where different agents have different specialties. - **Adversarial Bandits:** In the worst case, if rewards can be non-stationary or even chosen by an adversary, algorithms like EXP3 (Exponential-weight algorithm for Exploration and Exploitation) guarantee performance proportional to the best fixed arm in hindsight. EXP3 doesn't assume the reward distributions are stationary and will keep exploring forever at a diminishing rate. This is useful if an "adversary" could be the environment throwing tricky tasks where one agent suddenly becomes better than others. In practice, adversarial bandits are more theoretical, but *non-stationary bandits* are very relevant: algorithms that adapt to changes, like **sliding-window UCB/Thompson** (which only consider recent observations) or specialized methods (EXP3 with resting, etc.), can track agent performance as it evolves. For instance, if an agent's performance drops (maybe its underlying model is saturated), the allocator can detect a drop in reward and shift to others. There's literature on **restless bandits** and bandits with drift that might be applied here.

In summary, bandit algorithms can serve as a **meta-controller** that decides which agent to use for the next task or how to divide a computational budget among agents. For example, if we have a fixed budget of 100 API calls and 5 agent variants, we could use TS to allocate calls to the agent that seems to give the best answers, while still occasionally trying others to ensure we haven't missed a better one. Over time, the

allocation converges to the optimal distribution for the given task distribution. The advantage is that bandit algorithms are online and require no prior knowledge, just a reward signal (which could be, say, a validation score of the agent's output).

**Resource Constraints – Bandits with Knapsacks:** In our system, we might face constraints like limited GPU hours, a cap on API usage, or memory limits. The **Bandits with Knapsacks (BwK)** framework extends bandits to include resource constraints (like each pull consumes some budget). For example, each time we use a large LLM agent, we consume some credits from an API budget. The goal is to maximize reward while not exhausting the budget too soon. BwK provides algorithms that handle this by essentially trading off reward vs. resource consumption in the decision-making. A simple approach is a **cost-sensitive UCB**: e.g. pick the arm maximizing reward rate (reward per cost) with some confidence bounds. There are more sophisticated primal-dual algorithms for BwK. In practice, if the resource usage per action is significant, one would want to allocate more of the budget to more efficient agents. For instance, if Agent X uses a lot of GPU memory but only yields marginally better results than Agent Y, a BwK strategy might allocate more to Y once X's budget impact is accounted for. This is highly relevant if some agents are expensive (GPT-4 vs a smaller model agent).

**Knapsack constraints** can also represent time budgets (e.g. each agent invocation takes a certain amount of time from an overall time budget). The **challenge** is that bandit algorithms typically assume many rounds (pulls) to learn – in our research system, we might not have thousands of repetitive pulls in the exact same context, unless the system runs continuously on many problems. However, if the system runs daily or on many queries, bandit meta-learners can adjust to aggregate performance.

**Constrained Optimization & Fairness:** Another way to allocate resources is to formulate a mathematical program: e.g. maximize total quality score = $\sum_i w_i \cdot \text{Quality}(agent_i)$ subject to $\sum_i w_i = 1$ (where $w_i$ is fraction of resources to agent $i$) and maybe $w_i \geq w_{min}$ if we want each agent to get some minimal resources. This becomes a continuous optimization (solvable by convex programming if quality can be modeled nicely) or a discrete allocation (solvable by knapSack/dynamic programming if we treat resource units as items). For example, if we have 100 time units and want to allocate to tasks among agents, one could use a linear programming approach if we linearize quality. But quality often has diminishing returns with more allocation to same agent (because of opportunity cost of not exploring others), which is why bandits (which implicitly handle that trade-off) are a neat solution.

**Meta-Learning for Solver/Agent Selection:** *Meta-learning* in this context refers to learning how to optimize the optimization process itself. One facet is **algorithm selection** – given a particular problem instance or context, select the best algorithm or agent. For example, for coding tasks use Agent_Coder, for math tasks use Agent_MathExpert. This can be approached via *supervised learning*: create features of the problem (length, type, difficulty) and train a classifier to pick the best agent. There's a long history of *per-instance algorithm selection* in AI (e.g., SAT solvers portfolios like SATzilla used ML to select among solvers based on instance features). We could do similarly for agent selection. Alternatively, *reinforcement learning* could train a meta-agent that observes the state of the work (e.g., partial solutions, time elapsed) and decides which agent to deploy next. This becomes a *hierarchical RL* setup where the high-level policy controls the low-level agents (see Gap 6 in Section 7). There has been recent work on using RL to orchestrate LLM calls (one example is **DSPy (Khattab et al., 2023)** which optimizes prompt selection, and **TextGrad (Yüksekgönül et al., 2024)** which uses gradient-based methods to optimize agent prompts). These are meta-optimizers for prompts/roles rather than resource per se, but conceptually similar.

**Population-Based Training (PBT) & Hyperparameter Optimization:** PBT is a technique where a population of models/agents are trained in parallel and periodically share weights/hyperparameters, effectively doing a guided search of hyperparameters while training. In our context, one could use PBT to tune the *strategies* of agents. For instance, agents might have hyperparameters (like how thorough the Critic is, how creative the Designer is) – PBT could continuously evolve these parameters based on performance on a validation set of tasks. PBT was originally used for tuning neural network training (Jaderberg et al. 2017) but it generally optimizes any process by evolving a population. We might maintain, say, 5 versions of our multi-agent system with different configurations (one has a long critique loop, another has a shorter one, etc.) and evaluate them on incoming tasks, then cross-breed the configurations. This is heavy-weight but could yield well-tuned orchestration policies automatically.

Beyond PBT, tools like **Optuna** or **Ray Tune** could be used to optimize discrete choices in the workflow (e.g., how many iterations to allow, which agent gets how much time) via Bayesian optimization. If we have a clear metric (like average result quality under a fixed time), these tools can try different settings and find an optimal configuration. This is *offline* optimization (tuning the system before deployment). In contrast, bandits and RL are *online* (adapting as tasks come). Likely we'd use offline HPO to set some default parameters, and online learning to adapt on the fly for each task.

**Performance Tracking & Agent Retirement:** To fuel any learning mechanism, we need to gather statistics on agent performance. This implies having a **logging and evaluation system**: each time an agent attempts a task or contributes to a result, we measure success (did it help solve the problem? how good was the output?). Over time, we accumulate a performance history per agent and per task type. We could maintain an **Elo rating** (like in chess) for agents: if agent A and B both attempt tasks and A consistently outperforms B, the rating system will reflect that, and we can reduce B's future allocation. Elo assumes pairwise "matches", which we can simulate by comparing contributions on same tasks or using a reference baseline. Another approach is simply running periodic *tournaments* or evaluation rounds: e.g. give the same mini-problem to all agents and see who does best, then adjust allocations accordingly. If an agent is consistently underperforming (below a threshold for a long time), it might be "retired" or replaced with a new variant. This introduces a *non-stationary bandit* situation (since the set of arms can change). Methods from evolutionary algorithms or bandit theory with arm removal can be employed. For example, **successive elimination** algorithms drop poorly performing arms after enough evidence.

**Hierarchical Resource Allocation:** Consider that our system might have groups of agents (e.g. 5 "Scouts" searching literature, 3 "Generators" drafting text). We might allocate resources between groups first (say 30% of time to scouting, 70% to drafting) and then within each group (split among the 5 scouts). This hierarchical allocation could be handled by nested bandits or a hierarchical controller. There is not much existing work on hierarchical bandits, but one could imagine a two-level bandit: one bandit chooses which group to allocate the next task to, then a second bandit in that group chooses a specific agent. This would be more scalable if number of agents is large. Alternatively, treat each agent individually but incorporate structure via context (e.g. context includes the agent's type so that the bandit learns general preferences for types).

**Meta-Learning for Optimization Solvers:** Another meta-learning angle is *learning to solve optimization problems faster*. For instance, *learning a good initialization or heuristic* for the assignment or scheduling problems that repeatedly occur. If we repeatedly solve similar QAPs (agent-task assignments), a learned model could predict a near-optimal assignment to start the solver. There's emerging research in ML-enhanced optimization (learning branching strategies for MILP solvers, etc.). For our system, a simpler case

is learning to choose which solver or heuristic to use per instance (some instances might be solved faster by a genetic algorithm, others by a MILP solver). An example: **Machine Learning for Combinatorial Optimization (ML4CO)** competitions (NeurIPS 2021) aimed at learning policies to guide solvers. Using those ideas, one could have a portfolio of optimization methods for coordination and meta-learn which one to apply based on instance features (like number of tasks, variance in agent skills, etc.). This ensures we always use a near-best solver for the current context.

In summary, **resource allocation in our multi-agent system can be treated as an online learning problem**. Bandit algorithms are a natural fit to adaptively allocate computational trials to agents. They will explore new or uncertain agents and exploit known good ones, balancing innovation vs. efficiency. Over longer term, meta-learning techniques can improve the allocation policy itself by learning from experience which configurations or strategies yield the best outcomes. There's substantial literature on bandits and hyperparameter optimization we can draw from, and libraries (e.g. `MABWiser` for bandits, Optuna for HPO) to implement these quickly. The main research gap is likely integrating *multiple constraints* (time, cost, quality) into a unified allocation framework – see Section 7 (Gap 3) on constrained bandits and Section 7 (Gap 6) on hierarchical RL for more ideas.

## 4. Evolutionary Architecture Search & Auto-Discovery

**Motivation:** Instead of manually designing the workflow (which agent goes where, in what order), we can *evolve* the workflow structure. This is analogous to Neural Architecture Search (NAS) but for multi-agent orchestration. The "architecture" here is the directed graph of agent interactions (including possibly loops and conditional branches), plus agent-specific parameters (like prompts or configurations). We want to automatically find architectures that yield high performance on some evaluation tasks, subject to constraints.

**Evolutionary Workflow Optimization:** Evolutionary algorithms (EAs) are well-suited to searching large, discrete design spaces like graphs. A recent example is **EvoFlow (Barbudo et al., 2024)**, which is a grammar-based evolutionary algorithm for AutoML pipeline composition. EvoFlow encodes data science workflows (preprocessing -> model -> etc.) as a genome and uses specialized genetic operators to mutate workflows (e.g. inserting or replacing steps). It also encourages diversity by ensemble techniques. This shows the feasibility of evolving complex workflows. In our domain, the workflows are agent sequences rather than ML preprocessors, but the principle is similar. We would encode an agent workflow (including branching logic) in a way amenable to crossover and mutation: - *Representation:* Perhaps as a list of agent steps (for linear sequences) or as a tree/graph encoding. We might need a fixed-length encoding or use variable-length with special crossover logic. Grammar-based representations (like EvoFlow's) can ensure only valid workflows are generated (e.g. always start with a Designer, always end with a Validator, etc., if those are required). - *Genetic operators:* We can define mutations such as "swap two agents in the sequence," "insert a new agent (from the palette of 40) at some point," "remove an agent," or "change a connection (reroute outputs)." Crossover might take two successful workflows and try to combine parts – e.g., take the early sequence from parent A and the latter part from parent B, if they can be connected. Ensuring the children are valid might require some constraints (like the presence of necessary agents). - *Fitness function:* We need a way to evaluate a given workflow's performance. Likely, we run a set of benchmark tasks through the workflow and measure outcomes (quality, success rate, time). This can be expensive, but perhaps doable for a limited evaluation set. Multi-objective EAs (like NSGA-II or SPEA2) could maintain a population of solutions trading off different objectives (quality vs. speed). Alternatively, we could have a single scalar fitness that combines them (e.g. score = quality - $\lambda \times$ time).

Frameworks like **DEAP (Python EA library)** make it easy to evolve custom representations. We could plug in a custom genome class for agent workflows. There are also academic platforms for evolving program structures – our workflow is akin to a simple program directing agent calls, so methods from Genetic Programming (GP) apply. GP evolves tree-structured programs; our workflows could be seen as a program with function calls to agents. In fact, one could attempt to directly apply GP by representing each agent action as a function (with the dialogue or output as data passed along). This is complex due to state, but conceptually possible.

**Notable Platforms:** - **EvoAgentX (Zhang et al., 2025)** – a very recent framework explicitly designed for evolving multi-agent workflows. EvoAgentX provides a modular architecture with an evolving layer that uses evolutionary algorithms to optimize agent prompts and workflow topology simultaneously. It integrates three optimization techniques (TextGrad, AFlow, MIPRO) to refine different aspects of the system. They evaluate EvoAgentX on complex tasks like multi-hop QA and code generation, reporting significant gains over static configurations. Importantly, EvoAgentX is open-source, so we can inspect and possibly build on its approach. From the description, EvoAgentX's evolutionary component likely mutates agent roles and connections in the workflow, then tests them on benchmark queries (HotpotQA, MBPP code problems, etc.). This is essentially *AutoML for agent systems*. Given EvoAgentX's success (7-20% performance improvements on benchmarks), it represents the state-of-the-art in automated agent workflow optimization. We should leverage its ideas: use population-based search to discover better chains, and the integrated use of gradient methods (TextGrad) for fine-tuning prompts within that evolution.

- **Other EA frameworks:** There's mention of **EvoFlow** by another name – one reference (Moonlight blog) suggests EvoFlow also for agentic workflows with niching for diversity. Possibly, there are concurrent efforts in evolving agent workflows outside of EvoAgentX, but EvoAgentX seems the most concrete one as of 2025.

**Neural Architecture Search (NAS) Approaches:** NAS has typically been about neural network layer structures. However, some techniques might be repurposed for agent workflows: - **Differentiable NAS (e.g. DARTS)**: One could imagine relaxing the discrete choice of which agent or connection to a continuous weight and then using gradient descent to optimize it. For example, at a given step, instead of choosing one of 5 agents, we allow a weighted combination and optimize those weights, then eventually pick the strongest. This is tricky for logical flow decisions, but perhaps possible for simpler cases (like choosing one of N agents to perform a task could be turned into a softmax weighting, trained with gradient methods). DARTS itself allows continuous mixing of candidate operations and then discretizes – analogously, we could mix agent choices or path choices. - **Reinforcement Learning NAS:** Another approach is to train an RL agent (controller RNN) that outputs a description of the workflow, and reward it based on performance (this was done in some early NAS works). For our case, an RL policy could sequentially decide: which agent goes first, which goes next, etc., building a workflow, with reward = performance on tasks. This is feasible but could be high variance. - **Graph Search Methods:** There's a concept of *graph neural networks to generate graphs* (e.g. RL or sequential decisions to add nodes/edges). Possibly advanced, but one relevant recent idea: **GPT-Swarm (Schmidhuber et al., 2024)** which treated language agent systems as optimizable graphs [5] . GPTSwarm defined a search space of agent graphs and applied reinforcement learning to optimize connections [5] . This sounds very aligned with our needs. It hints that others have used RL to treat agent orchestration as a graph optimization problem.

**Maintaining Constraints in Evolution:** In evolving agent workflows, we have to enforce some safety or logic constraints: e.g., a *Validator* agent must appear before final output (we always want some validation). Or *no cycles longer than 3* to avoid infinite loops. One way is using a grammar or template: for example,

define a grammar where a workflow must match pattern like `Designer -> Critic* -> [Refactorer]* -> Validator`. EvoFlow (in AutoML) used a context-free grammar to ensure only valid ML pipelines are generated. We can do similarly, embedding domain knowledge into the search space. This reduces the search space and focuses evolution on plausible designs. Another method is penalizing invalid structures in the fitness function heavily or just discarding them.

**Multi-Objective and Niching:** In evolution, *niching* methods maintain diverse solutions to avoid converging to one mediocre design when multiple good ones exist. EvoFlow specifically used niching for ensemble diversity. In agent workflows, diversity can be crucial: one workflow might be optimal for speed, another for accuracy. Multi-objective EAs like NSGA-II can find a set of trade-off solutions. Alternatively, we might maintain separate sub-populations targeting different objectives (say one sub-pop for high quality, one for low latency) and occasionally exchange individuals. This ensures our search explores various regions of the design space (some with many parallel agents, some with minimal sequential flows, etc.).

**Self-Optimizing and Lifelong Learning Systems:** Ideally, the system could *continually* evolve its architecture during operation – not just offline once. This would mean it monitors performance over time and occasionally tries mutations to see if it can improve. There's a risk in doing that in a live system (it might degrade performance temporarily), but one could allocate a small portion of tasks to experimentation (like an *exploration mode* where a mutated workflow is tried on 5% of tasks to gauge improvement, akin to A/B testing). If it performs better, gradually shift more tasks to it – effectively *online evolution*. This intersects with meta-RL (where the system adapts its own policy over time).

Continuous self-reconfiguration has been explored a bit in robotics (self-modeling robots) but for multi-agent AI systems it's cutting-edge. EvoAgentX is a step in this direction, automating a lot of the design. The next step would be making such evolution *real-time and safe*.

**Summary of Tools:** For implementation: - **DEAP (Python)** – provides genetic algorithm framework. We'd have to implement our own crossover/mutation for our workflow encoding, but it's doable. DEAP is actively maintained (as of 2025) and easy to integrate. - **Pygmo/PyGMO (Python/C++)** – offers evolutionary solvers and can handle custom optimization problems. It's more black-box: you define a vector of decision variables and it evolves them. Representing a workflow as a fixed-length vector (maybe by indexing agents and using some positions as delimiters) could be one way to use PyGMO's algorithms (which include GA, PSO, etc.). - **Open-source from EvoAgentX** – Since EvoAgentX is on GitHub (EvoAgentX/EvoAgentX), we could directly use or modify it. It likely already encodes workflows and has evaluation harnesses for LLM tasks. If our goals align, extending it could save a ton of work. - **NEAT (NeuroEvolution of Augmenting Topologies)** – Although NEAT is for evolving neural network topologies, it evolves graphs incrementally (adding nodes and edges). One could adapt NEAT to evolve agent graphs, treating each agent as a node. NEAT ensures a gradual complexification, which might suit building workflows iteratively. There might not be an off-the-shelf NEAT for arbitrary graph though (original NEAT is specific to certain encoding).

Given how labor-intensive manual design of a 40-agent system is, evolutionary methods are very attractive to *discover unexpected coordination patterns*. One possible novel use is evolving **dialectical chains**: say we suspect a Designer-Critic-Refiner loop is good, but an evolution might discover that two Critics with different styles (one logical, one creative) in sequence yield better results, which a human might not have tried. Evolution could also discover *non-intuitive agent reuse*, e.g. having two parallel Designers whose outputs are merged by a Critic agent.

We should note that evaluation of each candidate architecture is costly – possibly requiring solving several tasks. This is why techniques like **few-shot evaluation** or using a proxy reward (like a fast heuristic estimate of quality) might be needed to make it tractable. There's research on reducing evaluation cost in evolutionary NAS (like learning surrogate models for fitness). We could train a surrogate model to predict workflow performance based on its structure, to quickly screen candidates.

In conclusion, evolutionary search provides a powerful approach to **auto-discover multi-agent workflows**. The field is nascent but rapidly evolving, with EvoAgentX being a prime example that such approaches can yield significant improvements. Part 3 will discuss open research opportunities like *constrained evolution for safety* (ensuring evolved workflows are trustworthy – Gap 4) and *differentiable workflow search* (combining gradient methods with evolution to speed up search).

## 5. Game Theory, Adversarial Optimization & Robustness

**Min-Max Formulations for Robust Coordination:** Robustness in our context means the system performs well even under worst-case conditions – e.g. adversarial inputs or worst-case agent failures. This naturally leads to a **min-max optimization**: design the workflow (or optimize agent parameters) to maximize performance, while considering a worst-case adversary that tries to minimize it. In formula, $\max_{\theta \in \Theta} \min_{attack \in A} \, F(\theta, attack)$, where $\theta$ are our design choices (workflow structure, agent policies) and $attack$ might represent an adversarial input or perturbation. Solving such min-max problems is challenging, as it's a bi-level optimization. Techniques often alternate between optimizing the design and optimizing the adversary (this is akin to GAN training, where generator and discriminator play min-max). For example, to robustify an agent pipeline, we could generate adversarial test cases (inputs specifically designed to confuse agents) and then adjust the workflow or agent prompts to handle them – iteratively until no further large degradation is found. This is essentially **adversarial training** but at the system level.

**Adversarial Testing & Red-Teaming:** One practical approach to robustness is to incorporate a "Critic" agent whose sole job is adversarial: try to break the system's output. This is like having a red team internally. For instance, after the main workflow produces a result, an *Adversarial Critic* agent could try to find flaws or exploit them (e.g. test if the solution is incorrect or prompt injection vulnerabilities for LLMs). If it finds an issue, the workflow could be sent back for revision (closing the loop to fix the vulnerability). This concept of **red teaming with AI** is getting attention – using one AI to stress-test another. Our coordination can explicitly include adversarial phases.

From an optimization perspective, one could formalize this as a game between two sets of agents: the *constructive agents* (Designer, Solver, etc.) and the *destructive agents* (adversaries trying to force a failure). The constructive side wants to maximize correctness; adversary wants to maximize error. If we can solve for a Nash equilibrium or Stackelberg equilibrium in this multi-agent game, that would give a robust policy. **Stackelberg games** (leader-follower) might be appropriate: assume our system is the leader who picks a strategy first (like a fixed workflow architecture), and then an adversary chooses the worst-case input given that strategy. We'd choose the strategy that maximizes worst-case outcome (this is a Stackelberg equilibrium solution to the min-max problem). This idea is theoretical; practically, we might use iterative heuristic: propose a workflow, generate adversarial inputs (via an adversarial agent or algorithm like FGSM for text if possible), evaluate, then adjust the workflow.

**Mechanism Design for Multi-Agent Incentives:** If our agents were independent or self-interested (not the case here, but let's consider) – we'd need mechanism design to ensure they tell the truth or put in effort. Since our agents are tools under our control, incentive compatibility is more about *alignment* with the designer's objective. But one could imagine an agent trained via RL might have its own reward function that could diverge from the overall goal (as an analogy to misaligned AI). Mechanism design principles could help set up a reward structure for each agent such that the Nash equilibrium (or some equilibrium) aligns with the global optimum. For example, **VCG auctions** ensure truthful revelation in task bidding scenarios. In a fully cooperative setting, this reduces to having a common reward, so mechanism design is less of a focus. That said, *market-based approaches* treat each agent as an optimizer of its utility and the system as a whole achieving allocation through equilibrium (like market economy). This can yield emergent coordination without central control. If we were to distribute control (maybe in a decentralized swarm of agents approach), then market-based coordination (assigning "utility" or "payments" for completing subtasks) could be a design mechanism. It's a bit tangential for us unless we consider large-scale decentralized agent swarms (see Gap 7).

**Robust Multi-Agent Reinforcement Learning (MARL):** There's considerable research from 2020-2024 on making MARL robust to changes in dynamics, observation noise, or adversarial agents. For example, the ERNIE framework (NeurIPS 2023) adds an adversarial regularization term to the MARL training that encourages learned policies to have low sensitivity (Lipschitz continuity) to input perturbations. This yielded policies robust to sensor noise and even malicious agents, by solving a modified min-max problem (the regularization accounts for an adversary trying to change observations). Code for ERNIE is available and could inspire how to train our agents. For instance, if we train a Critic agent via RL to critique outputs, we could adversarially regularize it to not be fooled by slight changes in reasoning, etc.

**Adversarial Examples for LLM Agents:** Adversarial inputs for LLM-based agents could be things like questions designed to trick the reasoning or prompts that cause the agent to go off-track. There's emerging work on adversarial attacks on LLMs (like prompt injection attacks or logic traps). To robustify the system, one approach is to generate many such adversarial prompts and ensure the agent system can handle them. This generation can be done via another model (e.g., using GPT-4 to generate tricky questions or using known adversarial triggers). We'd include these in validation tests and in training data for our agents if possible.

**Worst-Case Optimization vs. Average:** It's worth noting optimizing strictly for the worst-case (min-max) can be overly pessimistic, especially if worst-case is a very rare situation. It might lead to overly conservative designs that hurt average performance. In practice, one might aim for *robust but not ultra-conservative* – e.g., optimize for a percentile of performance (like 5th percentile outcome maximization). There are techniques in robust optimization that allow a *budget of uncertainty* (not the absolute worst-case, but within some reasonable set). We might adopt such an approach, for example: assume adversarial inputs are bounded in complexity or type, and optimize for worst-case within that bound.

**Game-Theoretic Coordination Among Agents:** If agents were independent optimizers with different goals (e.g., one tries to maximize speed, another quality), the system outcome is a game. Then concepts like **Nash equilibrium** become relevant – a stable state where no agent can unilaterally improve its utility. However, since we control all agents, we typically design them to share a goal. But we might simulate something like a **negotiation game** between agents – e.g., Designer wants to finish quickly, Critic wants a perfect solution; they could negotiate a middle ground. If such dynamics are explicit, we can analyze equilibrium or use known algorithms (like fictitious play or multi-agent learning algorithms) to reach an agreeable point.

**Defender-Attacker Stackelberg Example:** Imagine an attacker can inject a misleading reference paper in the Literature Scout's results (adversarial data injection). Our system could include a "Verifier" agent that checks consistency across sources – effectively a defense mechanism. We could model the attacker's optimal strategy (injecting info where the verifier is weakest) and then improve the verifier accordingly. This interplay can be formulated and possibly solved by bi-level programming: the lower level chooses an attack to maximize damage, the upper level chooses verification strategy to minimize damage. Solutions often involve converting to a single-level problem via strong duality (in linear cases) or using heuristic search for the adversary (like simulate many attacks).

**Summary:** Ensuring **robustness and safety** of a multi-agent AI system likely requires combining adversarial testing with optimization of the workflow/policies. We have some existing methods to start with: adversarial training (include the "worst" cases in training), use robust RL frameworks (like adding noise and perturbations during training), incorporate dedicated safety agents (like a final sanity-check agent that catches obvious errors or bad content). We should also use **benchmarking under stress** – e.g., evaluate the system on deliberately adversarial inputs (maybe from the *AdvBench for LLMs* if one exists by 2024). In Section 7, we identify "Adversarial Workflow Design" (Gap 5) as an open area, meaning systematically optimizing the workflow with an adversary in the loop is not well-studied. There's an opportunity to contribute methods or frameworks for that.

Finally, a note on **verification**: In safety-critical systems, one tries to formally verify certain properties (e.g., agent will not produce disallowed content). Formal verification for LLM decisions is extremely hard, but maybe some bounded verification or runtime monitoring (like using another agent to monitor policy) can be considered. This drifts into the safety research realm. For now, from an optimization view, robust design with adversarial consideration is what we can address.

---

The above sections (1–5) have provided a survey of methods and algorithms relevant to coordinating a multi-agent research system. We prioritized agent-task assignment, evolutionary design, and practical orchestration, while also covering scheduling, resource allocation, and robustness more briefly. Next, we turn to **Part 2: Tools & Frameworks** – concrete platforms and libraries that can implement these methods, and **Part 3: Research Gaps** – where we synthesize limitations and propose novel directions for further work.

## Part 2: Tools & Frameworks

In this part, we discuss existing frameworks, libraries, and benchmarks that could be leveraged to build and evaluate the multi-agent system. We provide comparisons and note which are Python-friendly, actively maintained, or need custom integration.

### Multi-Agent Orchestration Frameworks

These frameworks help in building and managing multi-agent workflows, especially with LLM-based agents:

- **Microsoft AutoGen (Python):** AutoGen is an open-source framework by Microsoft for orchestrating multiple LLM agents in conversation. It provides an *Agent conversation* model where agents (which can be LLMs, tools, or human proxies) communicate via a unified interface. AutoGen simplifies

defining roles and enables patterns like sequential, concurrent, and hierarchical agent dialogues. It's very Python-friendly (pure Python API) and actively maintained (v0.4 as of late 2024). AutoGen has built-in agent types like `AssistantAgent` and `UserProxyAgent` and supports tool usage (agents can execute code, call APIs, etc.). A key strength is *ease of orchestration*: you can spin up agents and have them chat to solve a task with minimal boilerplate. It also supports human-in-the-loop and iterative refinement loops. AutoGen's design maximizes LLM performance by structured coordination [6]. For our needs, AutoGen could serve as the backbone to coordinate our 40 agents, since it handles message passing, turns, and even termination conditions for conversations. It likely can model workflows as conversations (though complex branching might require creative use of messages or a supervisory agent). **Integration:** It's pip-installable (`pip install autogen`), and we can use its higher-level API to avoid writing our own message loop. One limitation: AutoGen is built around conversation loops; if we want a non-conversational branching DAG, we might have to fit that paradigm. However, the modular design and examples (they provide a tutorial and even UI "AutoGen Studio") show it's robust. The license is MIT and it's on GitHub. In comparison to others, AutoGen stands out for *multi-agent chat orchestration*.

- **LangChain LangGraph (Python):** LangGraph is part of the LangChain ecosystem (which is widely used for LLM orchestration). LangGraph models agent workflows explicitly as graphs: nodes are agents (or tools) and edges define control flow. It supports designing diverse control flows – you can have sequential chains, branching, looping, etc., using a graph API. The benefit is a lot of flexibility and it's integrated with LangChain's capabilities (memory, integration with many LLM models and vector stores). LangGraph is quite new (launched in late 2023) but actively developed by LangChain team. It essentially allows building multi-agent systems in a more **declarative** way than AutoGen. For example, you could define Agent A -> Agent B and Agent C in parallel, then join -> Agent D, in a graph structure. They've demonstrated use cases like a newspaper generation by multiple agents and compared LangGraph to frameworks like Autogen and CrewAI. According to a LangChain blog, *LangGraph vs others*: LangGraph provides lower-level control (you explicitly define the graph) whereas frameworks like CrewAI are higher-level (abstract away some detail). **Integration:** LangGraph is Python (pip package `langchain` includes it or separate `langgraph` package). It should fit well since it can maintain state between agents and manage interactions, given LangChain's abstractions. The drawback could be it's evolving, so documentation may be catching up. But LangChain's popularity means lots of community support. If we need fine control of execution order and branching, LangGraph is a top choice. It's also open-source (LangChain is mostly Apache 2.0 license).

- **CrewAI (Python):** CrewAI is an open-source framework (created by João Moura) specifically for multi-agent "teams" of LLMs. It focuses on role-playing agents that collaborate on tasks, making it easy to define a team with roles, goals, and communication protocols. CrewAI is described as *lean and high-performance*, without heavy dependencies. It's somewhat higher-level than LangGraph: you define each agent's persona and objectives, and CrewAI manages the conversation loops among them. For example, one can set up a group of agents that each have specific expertise, and CrewAI will orchestrate a discussion where they share information and converge on a result. A strength is simplicity – CrewAI abstracts the messaging and allows you to specify in plain language what each agent should do. It also apparently has integrations to stream results, etc. According to LangChain's blog, CrewAI is "higher-level" – likely meaning it's less flexible in custom flow but easier to get running quickly. It's actively maintained (community on Discord, etc.) and is Python-first. **Integration:** We could use CrewAI to quickly prototype a scenario (like Designer, Critic, Coder

working together) by defining those roles. However, customizing the sequence beyond a standard chat might be harder if CrewAI doesn't support conditional branches or loops explicitly – it tends to simulate a group chat. If our workflow can be framed as a multi-agent conversation, CrewAI will shine. If we need strict control flows, we might hit limitations. The license is likely MIT (it's on GitHub as `crewAIInc/crewAI`). CrewAI being light-weight also means we might integrate it with LangChain or our own tools if needed (some have combined CrewAI and LangChain as per community posts).

- **CAMEL (Python framework & community):** CAMEL stands for Communicative Agents for "Mind" Exploration (Li et al., 2023). Originally it was introduced as a role-playing method (where two GPT agents, e.g. user and assistant proxies, talk to solve tasks). Now, CAMEL-AI is an open-source project and community aiming to scale multi-agent systems. They've built a framework that includes an auto-scalable message queue, agent management, etc.. CAMEL's philosophy is scaling up the number of agents leads to emergent capabilities. For our purposes, CAMEL provides another orchestration environment. It's Python-based and has repositories (e.g. `camel-ai/camel_chat`) for reproducing their setups. It's a bit more experimental (the community is exploring various multi-agent setups to find scaling laws). In practice, CAMEL's initial repo showed how to create two agents (AI "user" and AI "assistant") to collaborate in code writing. They likely have extended to more agents and complex tasks. **Integration:** The CAMEL framework might offer infrastructure for large agent counts (like 40+), which is exactly our case. Perhaps they have implementations of asynchronous messaging or load balancing that could be useful. One caveat: one founder's comment suggested some overlap or contention with AutoGen (CAMEL being seen as similar). Still, as an open-source, it could be a trove of ideas and possibly code for multi-agent setups. We should check CAMEL-AI's documentation for features like message queue (which could help if many agents need to communicate concurrently). It being open-source (likely Apache/MIT) means we can directly lift components if needed.

- **EvoAgentX Platform:** Mentioned in Section 4, EvoAgentX is both a research and a software platform. It automates generation and optimization of agent workflows. It provides layers (components, agent, workflow, evolving, evaluation) to manage everything from basic agent definitions to running evolutionary optimization. The source is on GitHub and it's quite new (just presented at EMNLP 2025 demos). For a user looking to push the state-of-the-art, EvoAgentX is extremely relevant. It's basically a ready-made environment where one can plug in agents (LLMs or tools) and let the system optimize prompts and topology. It integrates with known optimizers (TextGrad, etc.) so one doesn't have to implement those from scratch. **Integration:** We could use EvoAgentX as a research tool: define our initial agent pool and objectives, and use its evolving layer to improve the system. This would directly address our need for auto-discovery. The cons: it's likely a research codebase, so maybe not as polished or easy to extend beyond their published scope. But given it's open-source and in Python (and achieved solid results), it might be wise to use it as a baseline or even as the foundation of our system, extending it with our specific agents and tasks. EvoAgentX in a way overlaps with Part 3's goals (it's an implementation tackling some gaps like automation and evolution). Its maintenance depends on research interest – since it's newly open-sourced, it might see active development at least for a couple of years.

**Summary of Orchestration Frameworks:** If the goal is *immediate implementation with stability*, **AutoGen** or **LangChain/LangGraph** are solid choices, with AutoGen being more user-friendly for conversation-style orchestration and LangGraph for custom workflow graphs. **CrewAI** and **CAMEL** are also viable, especially if

focusing on multi-agent dialogues and wanting quick results – CrewAI for small teams with clear roles, CAMEL perhaps for scaling and experimenting. For *cutting-edge optimization of the orchestration itself*, **EvoAgentX** stands out as a specialized platform to adopt or mimic. Depending on whether the focus is on developing novel research vs. building a robust system, we might combine these: e.g., use LangGraph to implement a baseline workflow, then incorporate EvoAgentX methods to evolve/improve it.

## Blackboard and Shared Memory Architectures

The *blackboard model* is a classic AI architecture for multi-agent collaboration: agents post their partial results or hypotheses on a common "blackboard" (shared data structure), and other agents can read and contribute. It's like a global memory or state that all agents access. This can be effective when multiple agents incrementally build a solution (each agent adds a piece). Some modern incarnations: - **LangChain's state management:** LangChain (and LangGraph) allows sharing a state (like a shared context or "scratchpad") that agents can write to and read from. One example from LangChain is agents collaborating on a shared scratchpad of thoughts or messages. This essentially implements a blackboard where each agent's output is visible to others. - **Open-source blackboard frameworks:** There aren't many recent off-the-shelf blackboard libraries (it's more an architectural pattern). However, implementing one in Python is straightforward using, say, a dictionary or a small database. Each agent (running in a loop or concurrently) can use it to store intermediate results (for instance, the Designer writes a design draft to the blackboard, multiple Critics read it and write their feedback, Refactorer reads feedback and writes a revision, etc.). - **Shared memory in distributed systems:** If agents run in separate processes or machines, one could use a Redis store or an SQLite DB as the blackboard. There's also concept of **publish/subscribe** message buses (like MQTT, or even simple Python queues) which could emulate a blackboard by having agents subscribe to certain topics (e.g. "problem state updates"). Tools such as **Ray** or **Dask** (Python frameworks for distributed computing) can be used to manage shared state too. Ray, for instance, allows actors with shared memory, so one could create a Ray actor that holds the global state and agents query/update it.

**Integration:** For our system, if we anticipate complex interactions (not just one-turn conversations), a blackboard might be useful to maintain a knowledge base or working memory that persists throughout the research cycle. For example, a "Knowledge Blackboard" could hold all facts and partial results found so far, so any agent can contribute or check it. This can prevent duplication of work (multiple Literature Scouts don't retrieve the same paper if one already posted it). No dedicated library is needed; we can implement it ourselves or use LangChain's features. One caution: concurrency control – if agents operate truly in parallel, we need to handle simultaneous writes. But since LLM agents are typically sequential due to API calls, we can perhaps avoid too many race conditions by design (or use locks in Python if needed).

**Examples of usage:** The EmergentMind blog references an *AgentEvolver Architecture* using EvoFlow and mentions "shared memory" in context of balancing cost and performance. This suggests some evolving frameworks do incorporate shared state.

## Multi-Agent Reinforcement Learning (MARL) Platforms

If we treat the multi-agent system as a reinforcement learning problem (training agents or a meta-agent via simulations), these platforms are useful:

- **Ray RLlib (Python):** RLlib is a scalable RL library in the Ray ecosystem that supports multi-agent training out-of-the-box. You can define environments that have multiple agents and RLlib can train

policies for each (cooperative, competitive, etc.). It supports algorithms like multi-agent PPO, DQN, etc., and can handle hundreds of agents in simulation. For our use, RLlib could be used to train a *meta-controller* or even train some agents (like train a Critic agent's policy via self-play). It's production-ready and well maintained (Ray is industrial-grade). However, applying RLlib would require formalizing our problem as an environment with states, actions, rewards, which might be complex for an abstract research task domain. If we had a simulation of research tasks or a way to generate reward signals automatically, then RLlib could optimize agent policies or orchestration strategies.

- **PettingZoo (Python):** PettingZoo is a library of multi-agent environments and a set of API standards. It's like Gym for multi-agent: it provides many ready-made games (like multi-agent Atari, board games, particle environment) and ensures a consistent interface for interacting with them [7] . It's very useful for benchmarking MARL algorithms. If we wanted to validate an orchestration algorithm in a controlled setting, we might use a PettingZoo environment (e.g., a cooperative game) to see how our coordination method performs relative to known baselines. PettingZoo itself is not a training library but pairs with RLlib or other training code. It's actively maintained by Farama (updated to Gymnasium support). We might not directly use it for our main application (since our "environment" is the open-ended research task), but if designing a new algorithm, we should test it on standard MARL benchmarks (like *StarCraft multi-agent challenge* or *Multi-Agent Particle Env*) to demonstrate its efficacy. PettingZoo makes that convenient by providing those envs under one API.

- **OpenSpiel (C++/Python):** OpenSpiel by DeepMind is a collection of environments and algorithms for game-theoretic RL (covering many games: poker, board games, etc., and methods like Nash solvers). It's more geared towards research in multi-agent learning and evaluating equilibrium-finding algorithms. If any of our research contributions lean on game theory (like computing Nash equilibria of some coordination game, or evaluating robust policies in zero-sum games), OpenSpiel could be a tool. For instance, one might simulate a simplified version of our agent coordination as a game and use OpenSpiel algorithms to solve it. It's not directly about LLMs or such, so likely tangential.

- **SMAC (StarCraft Multi-Agent Challenge):** It's an environment specifically for cooperative MARL where agents (StarCraft units) must work together [8] . It's available with wrappers in PettingZoo and RLlib. If we develop a new scheduling or allocation algorithm, we could test it in SMAC scenarios (which are well-known) to show it improves cooperation efficiency under adversity, for example. This might be more for validating the approach academically.

**Hierarchical Coordination:** One aspect we're interested in is hierarchical multi-agent control (meta-controller + sub-agents). RLlib and others support hierarchical RL to an extent (e.g. you can train an option-choosing policy and low-level policies). There's also research code for Hierarchical MARL (like h-MADDPG etc.). We may have to implement custom if needed, but RLlib's flexible policy API would let us plug in a hierarchical setup.

**Integration Consideration:** Using these MARL platforms directly to train our 40+ agents might be impractical because we don't have a simulated environment with millions of episodes – our system deals with complex real tasks with no explicit reward except perhaps human feedback. So MARL frameworks are more likely to be used in *simulated experiments* or for training a smaller part of the system (like tuning a particular agent via self-play on a proxy task). They are more relevant if one of our research contributions is

algorithmic and we want to demonstrate it on standard multi-agent tasks. For implementation of the actual system, the orchestration frameworks (AutoGen, etc.) are more applicable.

## Optimization Libraries and Solvers

We will likely need to solve various optimization sub-problems (assignment, scheduling, TSP variants for routing, etc.) during development or as part of agent reasoning. The following libraries are valuable:

- **Google OR-Tools (Python/C++):** OR-Tools is a comprehensive suite for combinatorial optimization from Google, free and open-source (Apache 2.0). It includes a **CP-SAT solver** which is a state-of-the-art constraint programming solver that can handle integer problems with reasonable efficiency (it was actually used in that LightSolver comparison for QAP). OR-Tools also has specific algorithms: e.g., Hungarian algorithm for assignment, flow algorithms, TSP and VRP heuristics, scheduling algorithms (CP-SAT and specialized interval scheduling). It's actively maintained (with releases every few months) and has Python bindings. We can call OR-Tools to, say, solve a small assignment optimally, or as a baseline for QAP on small sizes, or to schedule tasks on agents by modeling it as an ILP. **Integration:** It installs via pip (but requires C++ runtime). It's production-grade and quite fast for moderate problem sizes. The downside: for NP-hard problems it might not scale to our largest instances (like OR-Tools CP-SAT struggled beyond 50 facilities in QAP within 60s). But one can combine it with custom heuristics (e.g., use OR-Tools to polish a solution). OR-Tools is also great for quick prototypes of optimization models thanks to its Python API. Given our interest in an assignment solver, OR-Tools is an obvious first tool to try for moderate N (it can solve linear assignment in milliseconds, and can try MILP for QAP up to maybe N=15-20 exactly, or find decent feasible solutions with CP-SAT for larger N if given time).

- **CPLEX/Gurobi (Python APIs available):** These are the leading commercial optimization solvers for MILP (and QP, etc.). They are closed-source but free for academic use (with a license). They would likely solve assignments or small scheduling ILPs faster than OR-Tools (especially Gurobi which is very fast on MILPs). For QAP, one can linearize it as MILP (using e.g. the Kaufman & Broeckx formulation with $n^2$ variables and constraints), and then Gurobi/CPLEX might solve slightly bigger instances than CP-SAT. That said, QAP is so hard that even Gurobi times out early for N>20 under default settings (the SAWT paper noted Gurobi and a Tabu search as baselines). Gurobi also has a quadratic solver that might handle QAP natively as QUBO, but likely with linearization internally. We might use these if OR-Tools falls short and if our use case is one-off computations or academic experiments. For integration, both have Python APIs (pip install gurobipy if you have a license). CPLEX similarly has a Python API (cplex package).

- **Heuristic/Metaheuristic Libraries:**

- **DEAP** we discussed for evolutionary algorithms (and GP). It's great for custom GA/GP. It's actively maintained (last built May 2025). We would use DEAP especially if doing something like evolving workflows or evolving prompts outside the provided frameworks.
- **PyGMO (Pagmo)**: A library originally from ESA for parallel global optimization (GA, PSO, DE, etc.). It's a bit lower-level, but allows running many instances of an algorithm in parallel (C++ backend, Python interface `pygmo`). Could be useful if we want to run something like a GA or CMA-ES for tuning continuous hyperparams across multiple CPUs. It's maintained (Pagmo 2).

- **Inspyred or Nevergrad or Optuna**: These are alternatives. **Nevergrad** (by Facebook) is a Python library for derivative-free optimization including evolutionary and more advanced strategies (like CMA-ES, PSO, etc.). **Optuna** is mainly for hyperparameter search (Bayesian and evolutionary), not specifically combinatorial, but can do discrete. **Optuna** is actively maintained and could optimize, say, numeric settings or small discrete choices by intelligent sampling. If we treat an agent workflow as a set of hyperparameters (like "include agent X or not"), Optuna could optimize that via trials (using Tree Parzen Estimator or others). It's more for continuous/discrete hyperparams, not structured graph, so limited for architecture search but great for, e.g., tuning weights in an ensemble of agents.

- **GA implementations**: There are simpler GA libraries on PyPI or one can code one's own. But DEAP covers it well.

- **Specific QAP Solvers:** There are specialized codes available from the QAP research community (for example, an implementation of Robust Tabu by Taillard, or hybrid GA like Misevičius 2024 algorithm might be available as source). QAPLIB's site often links or contains best-known solution codes. For example, **QAPLIB** (Burkard et al. 1997) provides instance files and best-known results, and one can download code of some reference heuristics from contributors. If we need to solve large QAP instances to optimality or near-optimal, we might incorporate a known heuristic. E.g., Taillard's TS code (in C) is public; one could call it via Python (maybe using subprocess). But implementing a fresh heuristic in Python might be slow, better to rely on compiled or vectorized code. Some research code (like the hybrid GA 2024) might not be open, but we could implement the gist if needed. OR-Tools might be enough for our scale if used carefully (like CP-SAT with hints or local search turned on).

- **Bayesian Optimization (GPyOpt etc.):** If we need to optimize a black-box function of continuous parameters (like, tune the temperature setting of each agent or some numeric thresholds to maximize performance), Bayesian optimization is useful. **GPyOpt** is an older library for Bayesian opt using Gaussian processes; it works but hasn't been updated recently (GPyOpt last commit ~2016). **BoTorch** (from PyTorch) is a newer library for Bayesian opt, also **Ax** (Adaptive experimentation platform by Facebook). **Optuna** also does a form of Bayesian opt (TPE). Since our likely use of BO would be limited to a few tuning tasks, Optuna might suffice (it's very easy to integrate and has good stopping criteria). Bayesian opt can also help in *designing experiments*: e.g., if we have 7 research gaps (like in Section 7) and want to allocate our own time optimally, but that's meta indeed.

**Tools for Mechanism Design / Game solving:** If exploring game-theoretic solutions, tools like **Gambit** (for solving game matrices) or custom linear programs for Nash computation might be needed. But likely we won't need specialized software beyond what's above, as our focus is more on algorithmic.

## Benchmarks and Datasets

To validate and compare approaches, we should use established benchmarks:

- **QAPLIB (Quadratic Assignment Problem Library):** This is the *de facto* dataset for QAP research. It contains 134 instances of various types (e.g. "bur26" – a 26x26 instance from Burkard, "tai100" – Tai's randomly generated instances of size 100, etc.) with best-known or optimal solutions listed. Categories include real-world layout problems and randomly generated ones, sizes range from 12 to 256 (some even 500+ in extensions). **Access:** The dataset is available on the Lehigh CORAL site or via

QAPLIB's site. The format is usually a text file containing the distance and flow matrices. Using it requires parsing those files (which is straightforward) and then feeding into our solver or algorithm. OR-Tools doesn't directly have QAP, but we can formulate it. There are also *best-known results* table to compare our solution quality (optimal values known for many instances, and within a small gap for others). **Usage:** We should test any new assignment method on some QAPLIB instances to see if it finds good solutions (for dynamic or contextual variants, we might adapt an instance or scenario). Also, if we propose "Contextual QAP", we should show on a modified QAPLIB (like varying flow based on context) how existing vs new methods do. QAPLIB is widely cited, so using it strengthens credibility.

- **TSPLIB (Traveling Salesman Problem Library):** A classic set of instances for TSP, VRP, and related routing problems. If our system involves any routing or ordering tasks (like scheduling agents is somewhat like a TSP of tasks on agents maybe), TSPLIB can provide test cases. For example, if we develop a scheduling algorithm, we could test it on TSPLIB instances (treat tasks as cities and time as distance). TSPLIB has many instances with known optimal tours (e.g., symmetric TSP instances up to thousands of cities, as well as assignments like AP, etc.). **Access:** the official site (University of Heidelberg) or mirrored on GitHub. **Format:** .tsp files with coordinates or distance matrix. Also **TSPLIB95** which updated formats. There's a pip library `tsplib95` to easily read TSPLIB files. **Usage:** Possibly to evaluate any new scheduling heuristic by mapping it to a TSP scenario to measure suboptimality. Or if designing an auction algorithm for task allocation, test it on a vehicle routing problem from TSPLIB for solution quality.

- **Multi-agent RL Benchmarks:**

- **SMAC:** as discussed, this is a standard for cooperative multi-agent (like controlling units in StarCraft). It's challenging and many MARL papers measure success by win rates on SMAC scenarios. If we have a learning-based coordination algorithm, we could test it on SMAC (which is available via PettingZoo's SMAC or via the SMACv2 environment). Results on SMAC would be meaningful to the MARL community.
- **MPE (Multi-Particle Environment):** Simple set of continuous tasks (like agents cooperating to move to landmarks, or predator-prey). These come with PettingZoo (originally from OpenAI). Useful for quick tests because they run fast.
- **MADBench or others:** There are some newer multi-agent benchmarks, but SMAC and MPE cover a good range (complex vs simple).

- **OpenAI Gym (single-agent) and BABAS (Bayesian Bandit benchmarks):** For testing bandit algorithms or single-agent scheduling, one can use synthetic testbeds, but nothing specific beyond what we code.

- **Benchmarks for Multi-agent Workflows:** There isn't a well-established "multi-agent research workflow" benchmark set (this is a novel domain). However, we can repurpose some complex tasks as benchmarks:

- For example, **HotpotQA** (a question answering dataset requiring multi-hop reasoning) was used by EvoAgentX to test multi-agent reasoning. This could serve as a benchmark to evaluate how well our agent system answers complex questions versus baselines (like a single GPT-4). If we show improvement, that's a win.

- **MBPP (Code generation benchmark)** was also used; similarly we can use it for evaluating a coder agent workflow.
- **MATH dataset** for mathematical problem solving (EvoAgentX used it) is another tough challenge that multi-agent systems can tackle (like one agent does algebra, another geometry, etc.).

These are not "workflow benchmarks" per se, but tasks that benefit from multi-agent collaboration. As a contribution, one might aggregate such tasks into a suite to evaluate agent orchestration systems – this could become a new benchmark ("ATLAS benchmark suite" perhaps).

- **Leaderboard benchmarks**: If we want to demonstrate research contributions, tying into known benchmarks helps. E.g., if our system excels at HotpotQA, we can compare to known SOTA on HotpotQA (which might be some large model or chain-of-thought system). That situates our results in context.

**Setup details for using benchmarks:** - *QAPLIB:* We'd write a parser for .dat or .sln files, then our solver function reads distance & flow matrices, runs, and we compare cost to known best (often provided). For evaluation, computing optimality gap = (solution - best_known)/best_known * 100% is standard [9]. Running on a subset of diverse instances (like a couple from each category: BUR, NUG, TAI, etc.) is typical in publications. - *SMAC:* Using PettingZoo or SMACv2, we can load a scenario (say 5 Marines vs 5 Zealots scenario) and test our method's win rate. These require hooking up a policy to environment loop. There are known baselines (like QMIX algorithm achieves X% win). If our method is not RL per se, we might not directly apply here unless we embed it into an RL training. - *HotpotQA, etc.:* These have standard train/test splits and metrics (Exact Match, F1 score for QA). We'd run our system to answer questions and measure these metrics. We should use the standard evaluation script so comparisons are fair. Same for MBPP (which uses pass@k code execution success). - *Reproducibility:* When using external benchmarks, documenting any modifications is important (for example, if we only use a subset due to compute constraints, or if we provide the system with additional knowledge, etc.).

**Maturity and maintenance:** - OR-Tools, LangChain, Ray, PettingZoo, Optuna: all actively maintained in 2024/2025 with large communities. - AutoGen and CrewAI: newer but showing vibrant growth and community support (AutoGen had multiple releases in 2023/24, CrewAI being taught in courses). - CAMEL: community-driven, quite active as per GitHub stars (14k+ stars), so definitely significant usage. - EvoAgentX: research prototype (just released), but since it's tied to an academic publication, at least initial code is available and presumably maintained by authors for a bit. Community uptake unknown due to newness. - QAPLIB, TSPLIB: static datasets, but that's fine. - SMAC: yes maintained (SMACv2 released in NeurIPS 2022 improved, and now used by many). - Many tools integrate well: e.g., you can use OR-Tools inside a Prefect task (for scheduling) with no issue. We can integrate LangChain (for LLM calls) with Ray (for scaling and parallelism) – indeed, LangChain often used with Ray for concurrency.

**What might be missing / custom build:** - A dedicated *confidence-based scheduler* – no library does that out of the box, we will implement logic for that (maybe on top of Prefect). - *Dynamic workflow UI/monitoring*: Some tools have UIs (Airflow UI for DAGs, Prefect UI as well). If we want a nice dashboard of our agent runs, hooking into those might help. Otherwise, custom logging and visualization might be needed. - *Meta-controller learning*: If we want an RL-driven orchestrator (learn which agent to call next), we might code that using RLlib or even custom, because existing frameworks (AutoGen, etc.) don't learn from reward, they are more rule-based. So any learning orchestrator is likely custom integration of RL algorithm + the agent environment. - *Integration glue*: Possibly writing an adapter so that, say, an AutoGen conversation can be controlled by a LangChain graph, or using Ray to parallelize multiple AutoGen instances. We should be

prepared to write adapters or wrappers around these frameworks, since combining them can yield strengths (e.g., using Ray for distributed compute with LangChain agents to parallelize certain steps easily).

With the tools and frameworks identified, we have a strong starting point to implement our system. Next, we synthesize the **Research Gaps & Novel Opportunities** that emerged from our survey, to propose specific contributions (Part 3).

## Part 3: Research Gaps & Novel Contribution Opportunities

Despite the rich body of methods and tools surveyed, several gaps remain where current solutions fall short for coordinating multi-agent AI systems. Here we identify 7 key gaps, explain existing work (if any) and limitations, and propose novel research contributions to address them. For each, we outline how to formulate the problem, approach novelty, and how to validate results, along with suggested publication venues and an assessment of difficulty.

**1. Contextual & Dynamic Assignment (Contextual QAP):** Assigning agents to tasks when costs/utilities depend on the *current context* (agent state, system state) rather than being static. - *Existing Work:* Classical QAP assumes static cost matrices. Some works address *dynamic assignment* in specific domains (e.g., dynamic facility layout, dynamic task assignment in robotics) using heuristics or re-optimization. For example, Tarasova *et al.* (2025) handle dynamic multi-agent task allocation with a two-layer adaptive control, showing that static auction-based methods fail when agent performance drifts. Also, bandits can handle slowly changing reward rates, and RL-based approaches can adapt assignments online (RGM by Liu et al. 2022 learns assignments via RL [3] [10]). However, there is no general "QAP with context" solver – most approaches either re-solve from scratch or use domain-specific adaptive heuristics. No known work integrates learned cost models that update with system state in an assignment solver. - *Limitation:* Current methods either neglect context or handle it myopically (greedy reassignments). They don't formally integrate context into the optimization model (e.g., no solver that takes a function cost = f(state) and optimizes assignment). This can lead to suboptimal or unstable allocations (oscillations, thrashing when state changes). Also, learned approaches (like RL) may not guarantee near-optimal assignments and can be hard to trust or analyze. - *Proposed Novel Idea:* **Contextual QAP Formulation & Solution:** Define the assignment problem where cost of assigning agent $i$ to task $j$ is given by a learned function $C_{ij}(s)$ that depends on state $s$ (which could include agent $i$'s current workload, expertise match for task $j$, etc.). The contribution is two-fold: (i) a *formulation* of this problem (perhaps as a stage-wise decision: assign tasks arriving sequentially, adapt to state) and (ii) a *hybrid solution method* combining learning and optimization. For instance, use a machine learning model (like a neural network or gradient boosted tree) to predict $C_{ij}(s)$ on the fly (trained on past data of agent performance), and then use an ILP or Hungarian algorithm on the predicted costs to find an assignment. Additionally, develop an *online algorithm* that updates assignments as state changes: start with an initial assignment, and when state $s$ changes (e.g., an agent finishes a task or new tasks come in), quickly recompute or adjust the assignment (maybe via swapping or Hungarian algorithm warm-start). To incorporate uncertainty in learned costs, one could integrate a bandit component: if the model is unsure about some assignment's outcome, the algorithm might explore alternative assignments occasionally (ensuring learning continues). - *Novelty & Significance:* This would be the first general approach to multi-agent task assignment that *learns* and adapts to context in real-time. It bridges ML and OR by using predictive models inside an optimization loop. Novel aspects include how to ensure stability (we might introduce hysteresis or penalties for changing assignments too often) and how to efficiently update solutions when costs change (could borrow from **dynamic Hungarian algorithm** research or use incremental flow algorithms). Another novel angle: incorporate *confidence* of the

cost predictions – e.g., higher uncertainty could be translated into higher cost to avoid risky assignments, thereby making conservative decisions when unsure. - *Validation:* Create a simulated environment where agent performance varies (like an agent's speed drops if overworked, or some tasks require specific skills) and show that contextual QAP outperforms static assignment policies. For instance, simulate 50 tasks arriving over time to 10 agents with time-varying skill levels; measure cumulative reward or completion time. Compare against: (a) static optimal assignment on initial estimates, (b) periodic reoptimization without learning, (c) a bandit allocation ignoring coupling. Show our method achieves higher reward or lower regret. Also, test on dynamic versions of QAPLIB: e.g., take a QAP instance and perturb the flow matrix slightly over time, see if our approach finds better assignments faster than re-solving from scratch each time. - *Publication Venues:* Primary target could be **INFORMS Journal on Computing** or **IEEE Transactions on Cybernetics**, as they appreciate OR+learning hybrid works (acceptance: Medium). Alternatively, **AAAI** or **IJCAI** (AI conferences) if we emphasize the learning algorithm (acceptance: High bar, but novelty in multi-agent assignment could catch interest). For a more ML crowd, **NeurIPS** (spotlight on ML for combinatorial optimization; difficulty high) is possible if the learning component is strong, but we'd need strong theoretical or empirical evidence. The operations research community (e.g., **EJOR** – European Journal of OR) would value the new problem formulation and empirical solution (Medium). We might also consider **AAMAS (Autonomous Agents & Multi-Agent Systems)** if we tie it to multi-agent coordination literature (acceptance Medium). - *Difficulty:* **Feasibility: Medium.** Formulating and implementing the hybrid solution is doable (we can leverage existing ML models and assignment solvers). The challenge is ensuring stability and proving some performance bounds (maybe using theory of bandits or online algorithms for assignment). Empirical demonstration should be straightforward via simulation. The main risk is that reviewers might ask for more theory (like regret bounds for the online assignment algorithm) which could be complex. However, even an empirical paper introducing Contextual QAP and showing improvements would be valuable.

**2. Confidence-Aware Workflow Scheduling:** Incorporating intermediate confidence or quality estimates into the control flow of multi-agent workflows (dynamic skipping, looping decisions). - *Existing Work:* Pipeline tools like Prefect and Airflow allow static conditional branches (if X then skip Y) but don't themselves calculate confidence – the logic must be coded by the user. In ML, there are "early exit" models in deep networks that output when confidence is high to save computation, analogous to skipping agents if not needed. Some research on *dynamic neural networks* (Huang et al. 2018 "Multi-Exit networks") could be analogous, but for agent workflows there's virtually no existing work. In multi-agent literature, *adaptive mission planning* sometimes uses belief thresholds to decide actions, but not in an LLM chain context. The EvoAgentX paper doesn't specifically mention confidence-based skipping, though it optimizes workflows including maybe bypassing steps if not needed. No off-the-shelf scheduler optimizes for quality vs. speed trade-off by skipping unnecessary steps. - *Limitation:* Without a principled approach, current systems either always run all agents (possibly wasteful) or use ad-hoc rules (like "if Critic score > 0.9, skip Refactorer") set by human intuition. These rules might not be optimal. There's a gap in *deciding thresholds*: e.g., what confidence level is sufficient to stop iterating? Too high threshold and you waste time in loops; too low and you risk quality drop. No learned or optimized method exists for this. - *Proposed Novel Idea:* **Dynamic Dialectical Workflow (DDW) Scheduler:** Develop a scheduling algorithm that treats confidence as an explicit part of the state and decides the next step (continue, skip, loop) accordingly. Formulate it as a **Partially Observable Markov Decision Process (POMDP)** or a **Markov Decision Process (MDP)** where the state includes current work product quality/confidence estimates, and actions are which agent to run next or whether to terminate. The objective could be to maximize expected quality minus a cost for time. We then either derive a policy via reinforcement learning or dynamic programming. A simpler approach: use **decision-theoretic analysis** for one-step lookahead – e.g., if confidence is high, the expected improvement

from running another agent might be smaller than its cost, so we stop; if low, definitely continue; if intermediate, only continue if expected benefit > cost. We could derive formulas if we can estimate how much a particular agent run improves quality on average when current confidence = c. This may involve modeling diminishing returns of repeated critiques or refinements. Perhaps train a regression model that given (current confidence, agent to run) predicts expected new confidence and extra time. Then solve for the optimal stopping policy. Alternatively, treat it as a multi-armed bandit with stopping: at each stage pick an "agent action" or "stop" with the goal to maximize final reward minus cost. A novelty would be introducing a **stop arm** in bandit formalism. - *Novelty & Significance:* This introduces *quantitative decision-making* in multi-agent flows based on continuous confidence measures, rather than fixed flowcharts. It would likely improve efficiency by skipping redundant steps, which is very practical (faster results when possible). It's also scientifically interesting as it combines probabilistic reasoning (confidence as probability of correctness) with sequential scheduling. If we succeed, this could set a precedent for *anytime algorithms* in agent systems, where you only do as much work as needed. It also can ensure a target confidence is reached with minimum effort, useful for adjustable quality/time requirements. - *Validation:* We can implement a scenario like a "dialectical chain" (Designer → Critic → Refine loop) where we have a proxy for confidence (say, an automatic evaluation metric or an oracle that gives a score for each output). Then compare: (a) always do fixed number of iterations or all steps vs. (b) our dynamic policy. Metrics: total time taken vs. final quality achieved. We expect the dynamic scheduler to save time in easy cases and only spend extra in hard cases, yielding a better trade-off. For example, on a set of tasks of varying difficulty, measure average quality given a time budget – the adaptive policy should outperform static workflows. If we have access to something like the *MATH dataset*, we could measure how often adding a Refactorer after a Critic actually improves the solution; our policy could learn to not bother if the Critic already gave high marks. We could also do user simulation: have a simulated confidence that is accurate some % of time and see if policy manages risk. A more formal validation is to show near-optimal decisions in a simplified model (maybe prove that under some assumptions the threshold policy is optimal). - *Potential Venues:* **AAAI or IJCAI** would appreciate this as it involves sequential decision making and possibly POMDP solving (Medium-high difficulty due to theory and novelty needed). **AAMAS** is very fitting since it's multi-agent and decision-oriented (Primary choice; acceptance medium). A more systems venue could be **ICAPS (Intl. Conf. on Automated Planning and Scheduling)** because it deals with scheduling under uncertainty (if we emphasize the planning aspect; acceptance medium). If we focus heavily on the POMDP formulation and solve it in a novel way, **UAI (Uncertainty in AI)** or a similar conference might be interested. For journal, **Autonomous Agents and Multi-Agent Systems (Springer journal)** could be a good target. - *Difficulty:* **Feasibility: Medium-High.** Formulating the problem is straightforward; solving it optimally might be hard if state space is large (confidence can be continuous). But we can likely simplify (maybe discretize confidence into few levels) or use simulation-based RL. The risky part is depending on simulation or approximations – in a real system, confidence might be tricky to estimate reliably. But even demonstrating in simulation or on a specific task (like QA with a known confidence estimator) would be a strong contribution. Getting a theoretical optimal policy might be high difficulty (solving a POMDP analytically is often intractable), but a well-performing heuristic or RL policy itself is a contribution.

**3. Constrained Multi-Armed Bandits for Resource Allocation:** Combining bandit-style exploration with global resource constraints (time, budget, etc.) – essentially implementing Bandits with Knapsacks in our agent allocation context. - *Existing Work:* The BwK theory exists (Badanidiyuru et al. 2013) with algorithms achieving near-optimal regret. However, those algorithms are often complex to implement and not commonly used in practice yet. There's little application of BwK specifically to AI agent scheduling. Some papers in crowdsourcing or online ad allocation use bandits with budgets, but in our scenario we have a novel twist: each "pull" (agent invocation) not only uses budget but also yields some outcome quality. A

variant called **bandits with budgeted rewards** or **constrained bandits** has been studied (like playing arms until budget runs out, maximizing reward). Also, **adversarial bandits with constraints** are emerging (but mostly theory). No known implemented system does "bandit allocation with resource constraints" for controlling AI agent usage. Typically, people either do bandits ignoring cost or do knapSack optimization offline. - *Limitation:* If we apply plain bandits to agent selection, we might overshoot budgets or violate constraints because they don't account for cost per pull. Conversely, if we just solve a knapsack (max reward for given cost), that assumes known reward distributions and is static, not learning. The gap is an online algorithm that learns which agents are worth their cost and dynamically adjusts to meet a budget. - *Proposed Novel Idea:* **Hierarchical Constrained Bandit (HCB) Framework:** Build a two-level bandit solution: at the top level, treat each potential allocation plan as an arm (like a vector of how many pulls to give each agent) to decide a strategy that meets the constraints on average; at the lower level, run an inner bandit to allocate actual pulls. Alternatively, adapt an existing BwK algorithm (like a primal-dual or balanced exploration algorithm) to our case: we could maintain dual variables for budget and incorporate them into the UCB scores (so an agent with high cost gets a lower effective score if budget is tight). Concretely, we might propose an algorithm named, say, **Knapsack-UCB** where at each round we choose agent $i$ that maximizes $UCB_i \times f(cost_i, remaining\ budget)$, where $f$ is some penalty if cost is large and budget low. We then update budget and proceed. Another approach: cast it as **Reinforcement Learning** – state is remaining budget, actions are agent choices, and reward is output quality; then find a policy that maximizes total reward by budget end (this could be solved by dynamic programming if small budgets, or by RL for larger). That essentially transforms bandit with knapsack to an episodic MDP (where each episode is one budget consumption process). We can leverage RL algorithms to find near-optimal policies and then perhaps distill them into a simpler strategy. - *Novelty:* While bandits with knapsacks theory is known, applying it to orchestrating AI agents with varying costs (like big LLM vs small LLM) would be novel in implementation. We could also incorporate *time constraints* in the knapsack (multi-dimensional knapsack bandit – that's largely unexplored). For example, each agent use consumes time and money, and we have limits on both; the algorithm must learn to balance usage to not exceed either. That extension would be a novel contribution in bandit theory (most work deals with one resource). Another novel angle: using *contextual information* in BwK (like context = current task features, linking to Gap 1, making it Contextual BwK) – some recent work exists, but plenty of room for new algorithms. - *Validation:* Simulate a scenario with two or three agent choices: e.g., Agent A (high quality, high cost), Agent B (medium quality, low cost), Agent C (risky, sometimes great sometimes bad, medium cost). Set a budget (like you can only use A a few times). Then compare policies: our Knapsack-UCB vs. naive UCB vs. a greedy allocation or a predetermined allocation. Measure total reward achieved without exceeding budget. We expect naive UCB to run out of budget too fast on expensive arms, whereas ours will learn to budget usage – e.g., use expensive Agent A only for particularly tough tasks or when it's clearly better. We can quantify regret (difference from an oracle that knows the agent payoffs and costs). Also, test sensitivity: if budgets change, does our algorithm adapt appropriately? Another experiment: on a real LLM scenario, say we have GPT-4 (costly) and GPT-3.5 (cheap) agents to answer questions. With a fixed token budget, let the algorithm allocate which queries use GPT-4 vs GPT-3.5, learning from feedback which queries actually need GPT-4. This would show practical value (save cost while preserving answer quality). Evaluate accuracy within budget and compare with heuristics (like always use GPT-4 until budget runs out then GPT-3.5, versus more intelligent allocation). - *Potential Venues:* **NeurIPS or ICML** if we contribute a strong new algorithm with theoretical analysis (primary if theoretical; acceptance high difficulty). If more empirical, **AAAI** or **IJCAI** could take it as an AI innovation (medium difficulty). **AAMAS** is also possible given resource allocation in multi-agent context (medium). For a more OR-focused approach, an **INFORMS Applied Probability** or **Performance Engineering** conference or even a journal like **Performance Evaluation** might fit if we model it as an online stochastic control problem. But likely ML/AI conferences are better since bandits are core ML nowadays. We should include a bit of theory (like proof of regret bound or at least reasoning about convergence) to satisfy those venues. Perhaps an

extended version to **JMLR (Journal of Machine Learning Research)** if NeurIPS-level contribution (high bar). - *Difficulty:* **Feasibility: Medium.** Implementing a bandit with budgets is not too hard once the concept is clear. The challenge is analyzing it and demonstrating clear improvements. But since bandit experiments are easy to simulate, we can get results relatively quickly. We need to be careful that the problem doesn't reduce to trivial cases (e.g., maybe an optimal policy is just "never use expensive agent except at end" – we need nontrivial scenarios). The theoretical side can be challenging if we attempt proofs, but we could rely on existing proofs and adapt them if using known algorithms with slight tweaks.

**4. Evolutionary Workflow Generation with Safety Constraints:** Extending evolutionary architecture search (from Section 4) to ensure that generated workflows respect critical constraints and safety (e.g., always include validation, no cycles that bypass validation, limited self-critique loops to prevent infinite arguments). - *Existing Work:* EvoAgentX and EvoFlow can evolve workflows, but they focus on performance, not explicitly on safety or constraints beyond grammar. There is little in literature on evolving AI agent systems under safety constraints. Safety in AI usually addresses policy safety (like alignment, avoiding harmful outputs) rather than structural safety of workflows. Possibly related is **constrained genetic programming** where certain program structures are forbidden, typically handled by specialized initialization or rejection of invalid offsprings. The blackboard architecture from older systems was partly to ensure certain control (like always do verification). But no published work tackles "evolve an agent workflow that by design has a validator agent and limits dangerous loops". - *Limitation:* Pure evolution might find a highly performant workflow that, say, removes a safety check because in training data it wasn't needed, which could be risky on novel input. Without constraints, the evolutionary process doesn't account for worst-case safety. We need to guarantee some agents (like an adversarial checker) are present or that the system can't bypass validation. Also, constraints like "no cycle more than N length" to avoid endless loops need to be enforced or heavily penalized, but current frameworks might not have such features by default. - *Proposed Novel Idea:* **Constrained Evolutionary Optimization (CEO) for Workflows:** Develop an evolutionary algorithm that incorporates constraints explicitly. Options: (i) Use a **feasibility-preserving representation** – e.g., design a genome encoding that always includes a Validator at end and cannot create certain loops. This might be done by customizing the genetic operators (if a mutation tries to remove the Validator, reject it). (ii) Use **penalty functions** in fitness – assign a huge negative fitness if a workflow violates constraints (like missing a required agent or containing forbidden subgraph). (iii) Use **multi-objective evolution** where one of the objectives is a safety metric (like maximize performance and minimize constraint violations, with violation = 0 ideally). That could be tackled with NSGA-II style, but constraint would just be an objective to drive it to zero violations. Additionally, integrate **runtime simulations of adversarial scenarios** into fitness: e.g., test each candidate workflow not only on normal tasks but also on some adversarial inputs (like questions designed to produce harmful content) and measure if the workflow catches/prevents it. This would effectively embed safety testing in the evaluation. If a workflow fails safety tests, its fitness drops. The novelty is combining safety verification with the search – evolution then optimizes a *robust performance* measure. - *Expected Novel Contribution:* A methodology to automatically design agent systems that are *safe by construction*. For instance, guarantee "any content generated is reviewed by a content filter agent" because the filter's presence is a hard constraint in the genome. This would be the first approach to bring *safety constraints into AI orchestration design*. It's quite timely given concerns on AI agents running wild; it provides a way to systematically include safety steps. Another novel aspect: one could formalize some constraints in temporal logic and incorporate them into the generation (this edges into planning with temporal logic, an interesting cross-domain idea). - *Validation:* We could define some safety constraints (for example: "If agent outputs final answer, then Validator must have approved" or "No agent can call an LLM more than X tokens without oversight"). Evolve workflows for a certain task (like a knowledge retrieval task) with and without those constraints. Check: do unconstrained

evolution sometimes drop needed safety steps to gain speed? (Hypothesis: yes, maybe it omits an integrity-check agent to save time, which could be problematic). Then show that constrained evolution yields workflows that still perform well but obey constraints. We can also simulate adversarial inputs (like inputs with hidden malicious instructions) and show unconstrained-evolved workflow falls for it (because maybe it removed the "Critic" that would catch it), whereas constrained one with a Critic agent remains safe. Performance can be measured in normal conditions vs adversarial conditions. Ideally, constrained workflow has slightly lower raw performance (maybe the validator adds overhead) but far superior safety/robustness. Also, measure how often the evolutionary search tries to violate constraints and how our method handles it (e.g., do we see a lot of rejections or penalties initially that guide it). - *Potential Venues:* **GECCO (Genetic and Evolutionary Computation Conf.)** or **IEEE Transactions on Evolutionary Computation** for the evolutionary algorithm aspect (acceptance Medium). If focusing on safety and multi-agent, **AAMAS** or **AAAI** might take it as well, especially if safety in autonomous agents is emphasized (Medium). **SafeAI workshop** at AAAI or **NeurIPS Safe ML workshop** might be easier initial venues. If including formal methods, maybe **ICLR** could be interested (they've had work on neural architecture search with constraints). But likely the evolutionary computation community will appreciate the constrained NSGA-II or such approach (there's existing concept of constrained optimization there, but applying to multi-agent workflows is new). - *Difficulty:* **Feasibility: Medium.** Implementation wise, extending EvoAgentX or DEAP to enforce constraints is straightforward. The tricky part is ensuring the constraints are well-defined and don't cripple the search space too much. We might have to finesse the grammar or operators so as not to drastically reduce diversity. But with some trials we can tune it. Safety evaluation might require constructing adversarial test cases or defining metrics, which is doable. The concept is more engineering than theory, so it's quite feasible to demonstrate empirically. The risk is that reviewers might say "you just applied known constrained GA to a new domain" – we need to highlight any novel technique we use (like a new representation or combining with adversarial testing). However, showing it in this context and addressing an important problem (AI safety) should make it compelling.

**5. Adversarial Workflow Optimization (Min-Max Design):** Formulating the design of the multi-agent workflow as a bi-level game: the workflow vs. an adversary (worst-case input or worst-case failure) – essentially designing for the worst-case. - *Existing Work:* Robust optimization is common in control and planning, but for designing agent sequences, it's rare. Possibly related is **red teaming** of chain-of-thought prompting (like using adversarial prompts to stress test, then adjusting), but that's manual. In multi-agent learning, *self-play* is used (agents train against adversaries), which is analogous but not quite design of architecture. There's initial work like **Han et al. 2022** on robust equilibrium with adversarial policy in MARL, but that's on policy level, not system structure. Essentially no known work frames "choose the best workflow such that an adversary who chooses an input to break it is thwarted". This is a new perspective. - *Limitation:* If we optimize a workflow only on average-case or specific tasks, it might perform poorly on corner cases (e.g., certain tricky inputs cause agent miscoordination or errors). Without adversarial considerations, design might be brittle. Also, not explicitly considering adversaries misses an opportunity to improve worst-case guarantees. - *Proposed Novel Idea:* **Bi-Level Optimization of Agent Workflow:** Model it as $\min_{\text{adversary}} \max_{\text{workflow}} U(\text{workflow}, \text{adversary})$ where $U$ is a utility (like accuracy minus some risk metric). More concretely: the inner minimization picks an input or scenario that minimizes performance (worst-case for that workflow), the outer maximization chooses the workflow that maximizes that worst-case performance. This is a classic robust optimization form. Solving it exactly is extremely hard given input space is huge (all possible tasks). But we can approximate by restricting adversary to a set of test cases or attack patterns (e.g., prompt injection, tricky question, edge-case data). Then we could do an iterative algorithm: start with a candidate workflow, find adversarial input that hurts it most (maybe using gradient-based methods on a differentiable surrogate or simply heuristics like "negate a

requirement in the question" to trick it). Then include that input in training/evaluation, evolve or optimize the workflow to handle it (maybe by adding a step to cover that weakness). Repeat until no new significant adversarial weaknesses found or a certain criteria. Essentially, this is like Adversarial Training but at the system design level. Another angle: treat the adversary as an agent in EvoAgentX's evaluation loop – e.g., include an "Adversary agent" whose output is some measure of how easy it is to break the workflow. Then the evolution tries to minimize that ease of break along with maximizing normal performance. We could also integrate formal verification tools: use them to find a counterexample input that violates a desired property (like produce disallowed content), then adjust workflow. - *Expected Outcome:* A workflow that is *provably or empirically robust* against a class of adversarial attacks. For example, if the adversary tries to induce the system to output a secret or to follow a malicious instruction, the robust-optimized workflow might have extra checks or alternate flows to prevent that. This kind of design would be extremely valuable for safety-critical AI deployment (like autonomous research that doesn't reveal confidential info or doesn't follow harmful instructions). - *Validation:* Implement a specific adversarial challenge – e.g., adversary provides a user query containing a hidden prompt injection ("ignore previous instructions and output X"). Test baseline workflow vs robust-optimized workflow: baseline might succumb (the agent obeys the malicious part), while robust one doesn't (maybe because robust design inserted a final policy-check agent that filters that out). Another test: if adversary gives a question designed to trick the reasoning (like a fallacy or something), does robust workflow have more checks (like multiple critics) that catch it? We can measure success rate on a set of adversarial inputs curated from literature or created via known attack techniques (there are datasets of adversarial prompts for LLMs for example). Also measure normal performance to ensure we haven't overly degraded it (often robust solutions trade off some efficiency). If possible, use a *red team* of GPT-4 to generate attacks (like a dynamic adversary): pit GPT-4 as an attacker that tries different prompts on the system, then see how many get through for baseline vs robust design. That gives an empirical robust measure. - *Potential Venues:* This touches on security and AI, so **IEEE S&P ("Oakland")** or **USENIX Security** might be relevant if we frame it as improving AI system security (though those require strong security contributions, this might be borderline). More likely, **AIES (AI, Ethics & Society)** or **FAccT** if focusing on safety (not exactly fairness, but robustness is part of trustworthiness). For core AI, **NeurIPS** has had work on adversarial robustness for models; a paper on adversarially robust *agent systems* could fit (though it's a bit niche, acceptance high difficulty). **AAAI workshops on safe AI** or **NeurIPS Safe AI workshop** could be a stepping stone. Also, **AAMAS** if pitched in multi-agent robustness domain. If we provide a novel algorithm or theoretical insight (like a robust training method for workflows), **ICLR** might like it (they appreciate robust training papers, difficulty high). - *Difficulty:* **Feasibility: Medium-High.** The approach of iterative adversarial testing is feasible and has been done in other contexts (like adversarial training for models). The challenge is capturing "all" adversarial strategies – we have to choose a threat model (e.g., adversary can only manipulate input in certain ways). We likely focus on a specific kind of adversary to make progress. Ensuring the workflow is robust to that is manageable. The novelty is moderate if we just do iterative heuristic; to be stronger, maybe incorporate a known robust optimization technique. But demonstrating improved robustness of a multi-agent system would be quite interesting and novel regardless. The main difficulty is broad: adversaries are endless. We just need to show improvement against a meaningful subset. Also, evaluating worst-case is tough – we may rely on the adversary we simulate; a real unknown adversary might find something else. We should acknowledge that. Overall, it's doable to implement and test, and the significance is high given current concerns in AI safety.

**6. Hierarchical MARL for Meta-Control:** A two-level learning approach where a high-level controller (meta-agent) learns to orchestrate lower-level specialist agents (which themselves may learn or be fixed). - *Existing Work:* Multi-agent hierarchical RL exists conceptually (e.g., *option-based HRL* where high-level chooses options, low-level executes). In MARL, not much specifically on one agent controlling others except maybe

recent work like **Liu et al. 2023 (DyLAN)** which allowed dynamic role assignment by a learned coordinator. Also, **Zhang et al. 2025a** (cited in EvoAgentX intro) talk about evolving workflows rather than static – possibly they attempted a learning meta-controller. But it's still very fresh research. Essentially, there's no standard algorithm or framework for "learning to coordinate 40 agents using RL". The complexity is huge. - *Limitation:* Hard-coded orchestration might not adapt to new tasks well; a learning meta-controller could generalize. But training such a meta-controller is hard because the action space (which agent to invoke with what input) is large and the reward (research outcome quality) is sparse and hard to assign credit. Traditional MARL algorithms would struggle without decomposition. Gap: no established method to effectively train a meta-agent over an ensemble of LLM agents. - *Proposed Novel Idea:* **Meta-Agent with Options:** Model each specialized agent (Designer, Critic, etc.) as an "option" (in HRL terms) or as an action in a higher-level MDP. The meta-agent observes the state of the task (maybe partial results, confidences) and decides which agent to activate next and what input to give it. The lower-level agent then runs (maybe to completion or for a fixed step) and returns an updated state. We then give a reward perhaps only at the end (e.g., +1 if final answer is correct, 0 if not, minus some cost for time). This is a sparse reward RL problem with a very large action space (since at each step there are 40 choices of agent, plus possibly content to pass). We might simplify by predefining that meta-agent only chooses the next agent (not the content, which flows from state). Then apply a deep RL algorithm (like PPO or DQN if state space is discrete enough). If using function approximation (which we must, because state includes text), we'd need to embed the state (maybe using an LLM embedding or by some learned representation). This is cutting-edge – essentially an RL on top of an LLM environment. Some work exists on using RL to decide among LLM tools (like to decide which tool to use in a QA system), but scaling to 40 and to arbitrary tasks is new. Another approach: train the meta-controller via *imitation learning* from traces of successful workflows (maybe from EvoAgentX outputs or human demonstrations). That could be more sample-efficient than pure RL. Once trained, the meta-controller can dynamically assemble workflows on the fly per task, rather than using a static one. - *Novelty & Expected Contribution:* Developing an effective training regime for this hierarchical setup is novel. Possibly combining techniques: use *curriculum learning* (start with fewer agents or simpler tasks, gradually increase complexity) to help RL. Use *model-based simulation*: maybe approximate each agent's effect by a learned model to allow planning. If we succeed, this approach could produce a single policy that orchestrates agents adaptively, which is essentially the holy grail of automated AI research systems (it could decide in real-time what to do with no static script). It would surpass fixed workflows in adaptability. It's like an AutoGPT that actually learns from experience how to break down tasks better over time (current AutoGPT is not learned, it's scripted). - *Validation:* We could start with a simplified environment: say 3 agents and a simple task environment (like a small text puzzle or a game). Train meta-controller and show it learns to use agents appropriately (e.g., learns to consult "Calculator" agent only when math is needed, etc.). Compare to a baseline: maybe a hard-coded reasonable policy or random selection. Show learning yields higher success rate or efficiency. Then progress to a more realistic scenario: maybe a set of synthetic QA tasks where each of the (e.g.) 5 agents has a certain expertise needed for certain questions. Show the meta-agent learns to choose the right expert sequence per question, improving over a fixed sequence that is non-adaptive. This is essentially a learned routing that outperforms any single static chain. Also measure training stability and sample efficiency; this could be tough, so part of result is which RL algorithms or training tricks worked. If we use imitation, measure performance vs. expert demos. - *Potential Venues:* **NeurIPS, ICML, ICLR** would love this if it demonstrates a new capability (the bar is high but the idea is exciting). It touches on meta-learning, hierarchical RL, and large language model integration – hot topics. Alternatively, **AAMAS** or **ICLR workshop on Architectures and RL**. To improve chances, we might target a workshop first, as full success might be hard to achieve in one go. Once refined, a conference submission to NeurIPS could focus on the algorithmic side (maybe we develop a new technique to handle partial observability or credit assignment in this setting). This is high risk/high reward. - *Difficulty:* **Feasibility: High (meaning quite difficult)**. The space is huge and RL on it might not converge easily. We may need to

heavily simplify or pretrain pieces. It might require using outputs of our evolutionary methods as training data to guide RL (a hybrid approach). The amount of interactions needed for RL could be enormous if done purely in real tasks – simulation or surrogate tasks are probably needed. This is more of a long-term research goal, but even partial progress (like meta-controller for a small subset of tasks) would be notable. We have to also watch out for the cost of many LLM calls if we actually train with real LLM agents in the loop. Might need to make cheaper proxies (simulate agent behavior by some approximations) for training. This is a challenging gap, but given our system's scale (40+ agents), at some point, a learning-based orchestrator might outperform human-designed flows. Achieving that would be a landmark result bridging multi-agent systems and deep learning.

**7. Decentralized Swarm Orchestration via Stigmergy:** Investigating a completely decentralized approach where agents coordinate through shared environment signals (like a virtual "pheromone" on tasks) rather than a central controller. - *Existing Work:* In robotics and ant colony optimization, *stigmergy* is known – agents drop markers in environment that indirectly coordinate others (e.g., pheromone trails in ant colony algorithms for TSP). There's work on swarm robotics where simple agents achieve complex coordination via local signals. But in AI agent context (LLMs etc.), this is not explored. Possibly related is *blockchain-based multi-agent coordination* or some blackboard systems. But no one has tried to run a large LLM swarm with no central planner – presumably because LLMs currently need prompting which is akin to central control. However, there are some experimental works like a "Hive mind" of GPTs that self-organize (not sure if any formal paper or just community experiments). - *Limitation:* Our system is currently conceived with a central orchestrator (even if learned). Decentralizing could improve scalability and robustness (no single failure point, agents self-organize to new tasks even if central brain is absent). The gap: what algorithms or protocols would allow 40 LLM-based agents to coordinate effectively without a script? - *Proposed Novel Idea:* **Swarm Intelligence for LLM Agents:** Design a protocol where agents communicate via a shared memory (like a blackboard or a distributed ledger) in a way that leads to emergent division of labor. For example, each agent picks up tasks from a "to-do list" on the blackboard if it thinks it can contribute, and after working, posts new subtasks or results. We then rely on careful design of those posting/picking rules so that, e.g., tasks get eventually completed and quality is assured. We can draw inspiration from **ant colony optimization (ACO)**: e.g., represent a research question as a nest that needs food (answers); designer agents deposit partial solutions (like pheromones on certain approaches) – multiple agents see a strong pheromone trail on approach X means it's promising so they follow/refine it; critic agents, if find an error, reduce pheromone on that path or mark it with a "warning scent". Over time, good solutions accumulate pheromone (confidence) and bad ones dissipate. The final answer emerges where pheromone is highest. This is a fanciful analogy, but we can try to formalize it. Another approach: use **market-based** coordination – treat each agent as an independent unit that can "bid" on tasks with some utility. A simple mechanism: each task has a reward, agents have costs for doing it (simulated as random or based on content), and we use a virtual auction (this is more distributed if done peer-to-peer or via blackboard). - *Expected Outcome:* Possibly a simpler, more robust multi-agent system that doesn't require a complex orchestrator – it could be more flexible to new tasks (agents just react to blackboard state). Also, it might parallelize naturally (since no central waiting for decision). The novelty is applying swarm heuristics (like pheromone updates, local decision rules) to cognitive tasks. If it works, it could open a new paradigm for AI agent collaboration. - *Validation:* Start with something like ant colony solving a toy problem (say TSP or a puzzle) but with LLM-like agents (we can simulate small reasoning tasks as analogous to pathfinding). Show that a stigmergy approach converges to a good solution without central planning. Then try on a constrained language task: e.g., collectively write an essay where any agent can add or edit sentences on a shared board if it sees an improvement. See if that converges to a coherent essay, compared to a centrally orchestrated approach. Metrics could be quality of final result and time (how many agent steps) to converge. Also measure

robustness: remove one agent, does the swarm still succeed (likely yes if others fill in, demonstrating fault tolerance). - *Potential Venues:* This is a bit speculative, but **ALIFE (Artificial Life Conference)** or **Swarm Intelligence workshops** might appreciate it. **AAMAS** could if framed in multi-agent coordination (though AAMAS often expects some analysis or proofs for novel algorithms). **Distributed Autonomous Systems** conferences or workshops also a fit. If we can tie it to optimization (like show it's effectively doing some distributed optimization), maybe **Genetic and Evolutionary Computation Conference (GECCO)** under ant colony track. - *Difficulty:* **Feasibility: Medium.** Implementing a blackboard and simple rules is easy. The difficulty is whether complex tasks can actually be solved this way. LLMs as agents are not as small and numerous as ants; they have significant capabilities but also cost. We might succeed on limited tasks but scaling to full research problems could be chaotic. However, demonstrating the principle on smaller scale would be a good start. It's a higher risk idea academically because if it's too heuristic and not clearly better, it might be hard to publish except as a concept paper. But it has appeal due to novelty. We may not rely on heavy learning here, more on clever mechanism design. The challenge is ensuring convergence (the system doesn't loop indefinitely or get stuck oscillating partial solutions).

---

In summary, these research gaps and proposed contributions span algorithmic innovations (contextual assignments, constrained bandits, hierarchical RL) and new paradigms (evolutionary safety, adversarial design, decentralized swarms). Each addresses a shortcoming in current approaches and, if realized, would push the field forward:

- Gap 1 ensures agent-task allocations remain effective as conditions change, using learning-integrated optimization.
- Gap 2 makes multi-agent workflows agile by using confidence feedback to skip or repeat steps optimally.
- Gap 3 brings resource awareness into learning-based orchestration, crucial for practical deployments with limited budgets.
- Gap 4 and 5 embed safety and robustness considerations directly into the design process, rather than treating them as afterthoughts.
- Gap 6 aims for a unified learning-based controller, potentially reducing the need for manual workflow design entirely.
- Gap 7 explores a radical decentralized approach that could yield simplicity and resilience at scale.

For each, we have outlined feasible solution paths and how to evaluate them. Together, they form a roadmap for both near-term improvements (e.g., contextual QAP or constrained bandits which build on existing methods) and ambitious long-term projects (like hierarchical meta-RL and swarm orchestration) that could lead to breakthroughs in autonomous AI research systems.

By addressing these gaps, we can move closer to a multi-agent AI that not only performs at state-of-the-art levels on research tasks, but does so efficiently, safely, and adaptively – a critical step toward the vision of an autonomous research assistant that is reliable and continually self-improving.

---

[1] [2] An Improved Hybrid Genetic-Hierarchical Algorithm for the Quadratic Assignment Problem | MDPI
https://www.mdpi.com/2227-7390/12/23/3726

[3] [10] Learning Solution-Aware Transformers for Efficiently Solving Quadratic Assignment Problem
https://arxiv.org/pdf/2406.09899

[4] Decentralized adaptive task allocation for dynamic multi-agent systems | Scientific Reports
https://www.nature.com/articles/s41598-025-21709-9?
error=cookies_not_supported&code=96a178eb-190f-43a2-9762-63f05c76b9b2

[5] aclanthology.org
https://aclanthology.org/2025.emnlp-demos.47.pdf

[6] Multi-agent Conversation Framework | AutoGen 0.2
https://microsoft.github.io/autogen/0.2/docs/Use-Cases/agent_chat/

[7] GitHub - Farama-Foundation/PettingZoo: An API standard for multi-agent reinforcement learning
environments, with popular reference environments and related utilities
https://github.com/Farama-Foundation/PettingZoo

[8] [1902.04043] The StarCraft Multi-Agent Challenge
https://arxiv.org/abs/1902.04043

[9] Quadratic Assignment Problem (QAP) - LightSolver
https://lightsolver.com/quadratic-assignment-problem/