

CoCap: Coordinated motion Capture for multi-actor scenes in outdoor environments

Aditya Rauniar¹, Micah Corah², and Sebastian Scherer¹

Abstract—Motion capture has become increasingly important, not only in computer animation but also in emerging fields like the virtual reality, bioinformatics, and humanoid training. Capturing outdoor environments offers extended horizon scenes but introduces challenges with occlusions and obstacles. Recent approaches using multi-drone systems to capture multiple actor scenes often fail to account for multi-view consistency and reasoning across cameras in cluttered environments. Coordinated motion Capture (CoCap), inspired by Conflict-Based Search (CBS), addresses this issue by coordinating view planning to ensure multi-view reasoning during conflicts. In scenarios with high occlusions and obstacles, where the likelihood of inter-robot collisions increases, CoCap demonstrates performance that approaches the ideal outcomes of unconstrained planning, outperforming existing sequential planning methods. Additionally, CoCap offers a single-robot view search approach for real-time applications in dense environments.

I. INTRODUCTION

Unmanned Aerial Vehicles (UAVs) equipped with cameras have garnered considerable attention for their wide range of applications, including surveillance, search and rescue, environmental monitoring, mapping, inspection, and 3D reconstruction [1–6].

Motivation: Recently, UAV teams have demonstrated effective view gathering for 4D pose reconstruction of single actors in outdoor environments [6]. This is particularly significant as traditional motion capture, typically confined to indoor studios, faces space limitations and incurs high costs for setting up the stage for stunts or other video shoots. Further research has extended this approach to multi-actor scenarios, where sequential view planning by a team of cameras ensures diverse and comprehensive coverage of a group of actors [7]. Moreover, the development of a robust multi-camera aerial system still requires study.

Challenges: Such motion capture applications in the wild present various challenges and requirements. First, since multi-view capture has been shown to improve 4D pose reconstruction for multiple actors [8–10], this demands a system of multiple camera-equipped drones capable of capturing diverse views of moving actors. Second, natural environments are filled with obstacles and occlusions, especially when cameras are maneuvered to obtain optimal coverage. This requires the

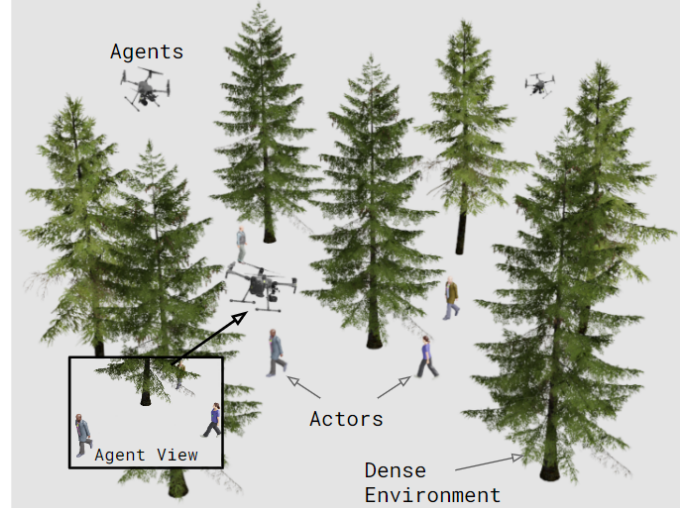


Fig. 1: **Coordinated View Planning:** Coverage optimization on dynamic actors with flying cameras in an occlusion-aware and obstacle-clustered environment where camera extrinsic positions across robots are negotiated.

drone system to not only manage obstacles and occlusions but also optimize for maximum coverage. Third, in these operational scenarios, as more camera-equipped drones are deployed, the system must also coordinate the views to avoid redundancy and prevent deadlocks in real-time.

CoCap: To address these challenges, we develop an efficient method for incorporating constraints into the system as more drones are deployed. Drawing inspiration from Conflict-Based Search (CBS) [11], we introduce constraints to the agents (drones) as conflicts arise. We define conflicts as instances where two agents enter a collision state, and we resolve this by applying constraints to one of the agents while maintaining a constraint tree. Our proposed approach of Coordinated motion Capture (CoCap) is tested in simulation, utilizing the reward structure employed by GreedyPerspectives [7]. Additionally, for real-time applications, we propose a single-agent view search method that prioritizes actor coverage and uses a heuristic to explore near-optimal viewpoints. More details on the simulation setup and reward structure can be found in our earlier research with GreedyPerspectives [7].

Contributions: Main contributions are as follows:

- A multi-camera view planning system that promotes diverse coverage from multiple viewpoints while avoiding camera placements focused on individual viewpoints.

¹A. Rauniar, and S. Scherer are with the Robotics Institute, School of Computer Science at Carnegie Mellon University, Pittsburgh, PA, USA {rauniar, basti}@cmu.edu

²M. Corah is with the Department of Computer Science at the Colorado School of Mines, Golden, CO, USA micah.corah@mines.edu

This work is supported by the National Science Foundation under Grant No. 2024173 and supported by Defense Science and Technology Agency Singapore contract DST000EC124000205.

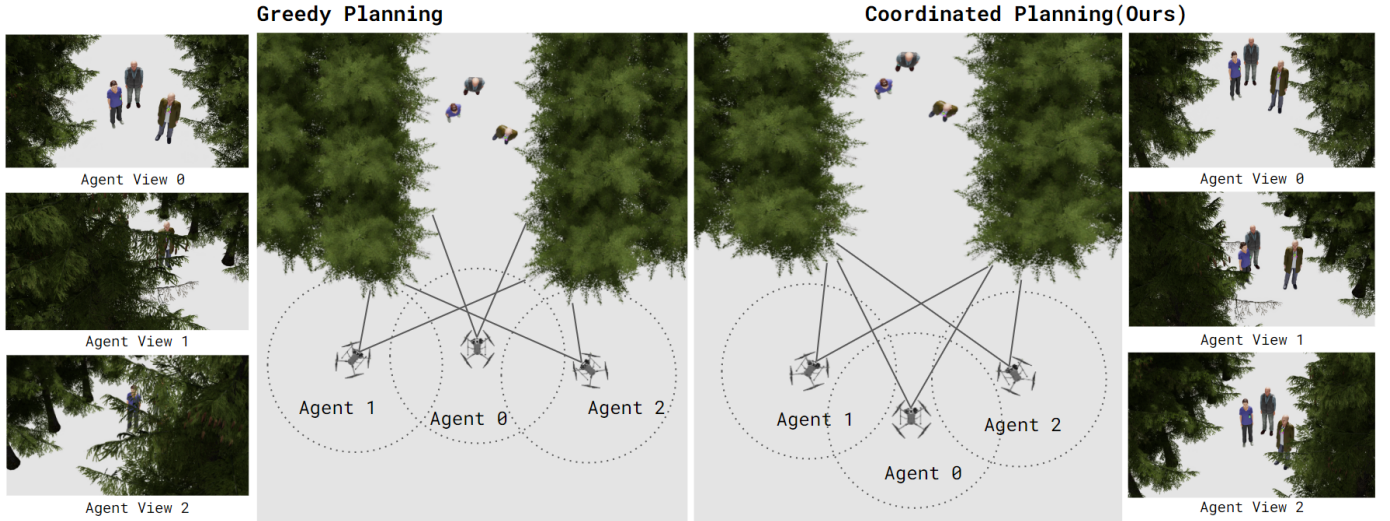


Fig. 2: **Sequential (Greedy) View Planning**: On the left, there is the sequential view planning of multiple camera positions, where there are egocentric behaviors across multiple viewpoints as seen in the three camera outputs on the left under greedy planning. **Coordinated planning**, on right: we propose a coordinated view planning approach where there is pixel-level negotiation amongst view positions to allow non-egocentric behaviors as seen in the three camera outputs on the right under coordinated planning.

- Scenarios of multi-actor obstacle clustered environments that are particularly challenging for multi-camera equipped drones.
- Heuristic-based single agent view planning framework for real-time applications.

II. RELATED WORK

Motion capture of actors in natural settings: Recently, there has been growing interest in developing motion capture systems using aerial cameras, particularly for outdoor environments. AirCapRL [12] introduces a deep reinforcement learning approach that learns an optimal policy for camera formations, focusing on achieving ideal viewing angles. However, this method struggles with reasoning about obstacles and occlusions in the environment and does not address scenarios involving multiple actors. Ho et al. [6] tackled the first limitation by developing a formation planner constrained to optimal views around an actor using a spherical grid, though it still does not account for multiple actors. One approach that addresses the multiple actor scenario is by Hughes et al. [13], who maximize Pixels-Per-Area (PPA) over multiple actors represented as polygonal cylinders. The GreedyPerspectives [7] further improves this method by addressing obstacles and occlusions. However, although these approaches can handle multiple actors, they struggle to maintain consistent visual alignment across multiple camera views in high-obstacle environments and are not optimized for real-time performance. CoCap addresses these challenges by introducing system constraints as conflicts arise and employing a heuristic-driven view planning method, guided by rewards from maximizing coverage across multiple actors.

Collision-Free Navigation for Multi-Robot Teams: Designing dense multi-robot systems presents significant challenges, particularly in conflict resolution. One widely used approach is prioritized planning, where robots are assigned a fixed priority sequence and planned sequentially based on their priority ID. This method has been employed in multi-actor view planning, such as in GreedyPerspectives [7]. However, prioritized planning can produce suboptimal in constrained and conflict-rich scenarios, as the preset order can cause later robots to encounter an increasing number of constraints, potentially resulting in deadlocks or planning failures. A more efficient approach to constraint development is Conflict-Based Search (CBS) [11], where constraints are dynamically introduced as conflicts arise during the expansion of a constraint tree. In multi-robot systems, conflicts typically occur when robots attempt to occupy the same position at the same timestep. Despite the success of CBS in navigation, its application in perception planning systems remains under-explored. CoCap extends the CBS approach to motion capture in multi-actor scenarios, specifically in outdoor environments with complex obstacles and occlusions.

III. PROBLEM FORMULATION

We consider a multi-robot system equipped with gimballed cameras that aims to monitor dynamic actors in an environment with obstacles. Each robot moves within a finite graph, selecting actions based on rewards tied to visual coverage. The following notation formalizes the setup, including robots, actors, and their observed faces, along with the workspace and reward structure.

Consider a set of robots $\mathcal{R} = \{1, \dots, N_r\}$, and a set of actors $\mathcal{A} = \{1, \dots, N_a\}$ each with a set of faces $\mathcal{F}_j =$

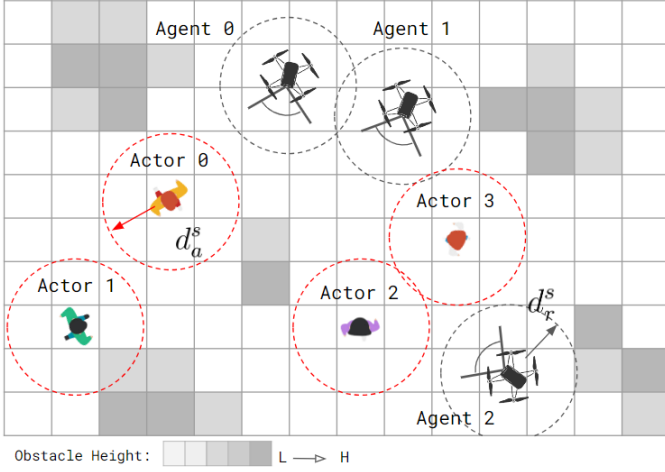


Fig. 3: Problem representation of the gimbaled camera (also formulated as a robot, \mathcal{R}) with projection matrix facilitating coverage on dynamic targets with occlusion and obstacles.

$\{1, \dots, N_{j,f}\}$ where $j \in \mathcal{A}$. Also, let all the sets of faces across all the actors be $\mathcal{F} = \{\mathcal{F}_1, \dots, \mathcal{F}_{N_a}\}$. Each of the faces, $\mathcal{F}_j : j \in N_a$, are associated with a set of pixel coverage by agents, $\mathcal{P}_x = \{\mathcal{P}_{x_{111}}, \dots, \mathcal{P}_{x_{ijk}}\}$, where $i \in N_r$, $j \in N_a$, and $k \in N_{j,f}$. All the robots move in a workspace represented as a finite graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, and share a global clock that start at $t = 0$. Vertex set \mathcal{V} represents the set of all the possible locations of i^{th} robot at time t as $v_t^i \in \mathcal{V}$ and the edge set \mathcal{E} as actions $e_t^i \in \mathcal{E}$ of i^{th} robot at time $t \in \{0, \dots, T\}$ where $i \in \mathcal{R}$. The reward for each edge is a M -dimensional nonnegative vector reward(e_t^i) $\in \mathbb{R}^+ \setminus \{0\}$.

Let $\xi^i(v_1^i, v_l^i)$ be a path that connects the vertices v_1^i and v_l^i via a sequence of vertices $(v_1^i, v_2^i, \dots, v_l^i) \in \mathcal{G}$. Let $g^i(\xi^i(v_1^i, v_l^i))$ denote the M -dimensional reward vector associated with the path $\xi^i(v_1^i, v_l^i)$ as the sum of all the reward vectors of all the edges present in the path, i.e., $g^i(\xi^i(v_1^i, v_l^i)) = \sum_{j=1, \dots, l-1} (v_j^i, v_{j+1}^i)$, in short, $g = J(\Xi)$.

Let $v_o^i, v_f^i \in \mathcal{V}$ represent the initial location and destination of robot i respectively. For simplicity, we denote a path from v_o^i to v_f^i for robot i as ξ^i . A joint path for all robots, denoted as $\xi = (\xi^1, \xi^2, \dots, \xi^{N_r})$, is referred to as a solution.

Objective: We employ submodular optimization techniques to synchronize the actions of a multi-robot team and to study the associated objective. The observation of each actor j is quantified in terms of the pixel density ($\frac{px}{m^2}$) obtained from a linear camera model's image. We define two functions, $\text{cov}(x_{i,t}, j, f) \rightarrow \mathbb{R}$, returning the pixel density for a specific actor's face observed from a robot's position, and $\text{covp}(j, f) \rightarrow \mathbb{R}$, returning the cumulative pixel density from all past observations. This enables us to express the incremental coverage gain:

$$\text{covm}(x_{i,t}) = \sum_{j \in T} \sum_{f \in \mathcal{F}_j} \sqrt{\text{cov}(x_{i,t}, j, f) + \text{covp}(j, f)} - \sqrt{\text{covp}(j, f)} \quad (1)$$

Algorithm 1: Coordinated View Planner

Data: X_{init} , \mathcal{T} trajectories, envHeightMap , agentMaxMotion

Result: U^\dagger

```

1 Initialization()
2 TREE  $\leftarrow$  GreedyPlanningWithoutConstraints()
3 Initialize  $U_{\text{seq}} \leftarrow \{\}$ 
4 while TREE not empty do
5    $P_k \leftarrow (\xi_k, \Omega_k, \mathcal{P}_{x_k}, g_k)$  // TREE.pop()
6   if no conflict detected in  $\xi_k$  then
7     return  $P_k$ 
8   end
9    $\Omega \leftarrow$  Split detected conflict
10  forall  $\omega_i \in \Omega$  do
11     $\Omega_{li} \leftarrow \Omega_k \cup \{\omega_i\}$ 
12     $\mathcal{P}_{x_l} \leftarrow \mathcal{P}_x \setminus \{\mathcal{P}_{x_{itk}}\}$ 
13     $\xi^{i*} \leftarrow \text{LowLevelViewSearch}(i, \Omega_{li}, \mathcal{P}_{x_l})$ 
14     $\mathcal{P}_{x_l}^* \leftarrow \mathcal{P}_{x_l} \cup \text{GetCoverage}(\xi^{i*})$ 
15     $\xi_l^* \leftarrow \xi_l \setminus \{\xi^i\} \cup \{\xi^{i*}\}$ 
16     $g_l \leftarrow \text{GetObjectiveValue}(\xi_l^*, \mathcal{P}_{x_l}^*)$ 
17     $P_l \leftarrow \{\xi_l^*, \Omega_{li}, \mathcal{P}_{x_l}^*, g_l\}$ 
18    TREE  $\leftarrow P_l$ 
19  end
20 end
21 return  $U_{\text{seq}}$ 

```

IV. BRIEF OVERVIEW ON CONFLICT-BASED SEARCH

Conflict-Based Search [11] is a two-level search that creates a binary tree to resolve conflicts in the agent paths on a high level and runs an optimal low-level search algorithm for a Multi-Agent Path Finding (MAPF) problem.

Conflict Resolution: Consider a pair of agents, i and j , each following their respective paths ξ^i and ξ^j . To detect conflicts between these paths, we utilize the function $\Psi(\xi^i, \xi^j)$. This function returns either an empty set if no conflict is present, or it provides details about the first conflict encountered along the paths. A conflict occurs at time t when agent i is at position v_t^i and agent j is at position v_t^j . This conflict is denoted by (i, j, v_t^i, v_t^j, t) . To prevent such conflicts, a corresponding constraint is added to the path of either agent i or agent j . This constraint is represented as $\omega^i = (i, u_a^i, u_b^i, t)$, where u_a^i and u_b^i are selected from the set of possible locations V . This constraint is associated with agent i and ensures that at time t , agent i follows the specified path, thus avoiding the potential conflict.

Given a set of constraints Ω , let $\Omega_i \subseteq \Omega$ represent the subset of all constraints in Ω that belong to agent i (i.e., $\Omega = \bigcup_{i \in N_r} \Omega_i$). A path ξ_i is consistent with respect to Ω if ξ_i satisfies every constraint in Ω_i . A joint path ξ is consistent with respect to Ω if every individual path $\xi_i \in \xi$ is consistent.

Two level Search: In the high level while creating the constraint tree with each node containing $(\xi, \Omega, g(\xi))$, each of them is in a priority queue. Initially, parent node $\mathcal{P}_{x_0} \in \text{TREE}$ is computed for all the agents such that the constraint set, $\Omega_o = \phi$, and $P_o = (\xi_0, \Omega_o, g_o)$ gets pushed to TREE. Let, the

total number of nodes created be denoted as N_G , and the total number of nodes expanded as $N_E : N_E \leq N_G$.

As each node, k , in the TREE gets explored based on the least g value, $P_k = (\xi_k, \Omega_k, g_k)$, $\Psi(\xi_k^i, \xi_k^j)$ gets called for all the sequence pairs $(i, j) \in N_r : (i \neq j)$. If there is no conflict detected, the solution, ξ , is found and the algorithm terminates. If there is a conflict detected, (i, j, v_t^i, v_t^j, t) , constraints to the agents are created in the following way. Each condition is created where constraint is added to either i or j , such that $\omega^i = (i, u_a^i = v_t^i, u_b^i = v_t^j, t)$ and $\omega^j = (j, u_a^j = v_t^j, u_b^j = v_t^i, t)$. Each of these constraints leads to a new node in the TREE, where $P_{li} = (\xi_{li}, \Omega_{li}, g_{li})$, and $P_{lj} = (\xi_{lj}, \Omega_{lj}, g_{lj})$, such that $l^i \in N_G$ where new set of constraints is added to i or j . Here, $\Omega_{li} = \Omega_k \cup \omega^i$, and $\Omega_{lj} = \Omega_k \cup \omega^j$. In each of these cases, for agent i , the algorithm updates the path ξ_k^i in ξ_{li} using the low-level search with a set of constraints Ω_{li} . A similar call is made for agent j leading to ξ_k^j . If the low-level search is unable to find a solution for any of these cases, respective P_{ln} is discarded.

CBS solves a single objective cost function, g , optimally by iterative expansion of a node in the constraint tree based on the least g val, resolving first agent-agent conflict (if exists) and creating a new node in the tree with new Ω , or returning the solution $P = (\xi, \Omega, g)$.

V. COORDINATED VIEW PLANNING

Our proposed planner Alg. 1 also works in a similar structure to CBS. With the following key differences:

- **Greedy planning:** Rather than initializing with solutions to single-agent path-finding problems we initialize with an approximate solution to the joint multi-agent view planning problem that relaxes constraints between robots.
- **Tree Node (P_k):** This node tracks pixel coverage by agents over actor faces, defined as $P = (\xi, \Omega, \mathcal{P}_x, g)$, where ξ is the joint robot path, Ω the path constraints, \mathcal{P}_x the total observed pixels, and g the node's reward.
- **Cost(g):** The cost, g , is the overall reward of the tree node as denoted by eq. (1).
- **Termination criteria:** Considering the computational expense, we return the first set of the solution P without any conflicts. If there is no valid solution, a failure is reported.
- **Replanning:** As a new conflict between any two agents (i, j) is detected using Ψ , the low level plans for agent i , with constraint set Ω_{li} and $\mathcal{P}_x \setminus \{\mathcal{P}_{x_{ijk}}\} : j \in N_a$, and $k \in N_{j,f}$. Similarly, for agent “ j ”.
- **Low-level search:** Replace A^* with a reward-guided solver using value iteration, focusing on coverage over actors and stable camera motion—key criteria for single-agent view planning.

A. Constraint Tree Formation

Initialization: During this state, with $(\Omega_o = \phi, \mathcal{P}_{x_o} = \phi)$, as we plan for a series of agents, $i \in \mathcal{R}$, we plan each agent greedily towards maximizing its reward using the single-agent view planner to produce the set of trajectories ξ_o . The root

node, $P_o = \{\xi_o, \Omega_o, \mathcal{P}_{x_o}, g_o\}$ gets added to TREE. As each agent, j , gets planned, the pixel context as passed such that $\mathcal{P}_{x_j} = \{\mathcal{P}_{x_{jkl}}\} : j < i, k \in \mathcal{T}, l \in N_{k,f}$.

Finding a solution: As a node is popped from the TREE, $P_k = \{\xi_k, \Omega_k, \mathcal{P}_{x_k}, g_k\}$, Ψ checks for conflicts amongst all the sequence pair of agents $(i, j) : i, j \in N_r$, and $i \neq j$. If there are no conflicts found, P is the solution, and the high-level search terminates. There are generally multiple conflict-free solutions to a set of valid trajectories maximizing g . However, we are interested in the first solution for computational performance and guide the tree expansion by the objective function $g = J(\Xi)$.

Resolving conflicts: If there exists a conflict between any two agents, $(i, j) : i \neq j$, then agent i is re-planned with constraint set $\Omega_{li} = \Omega_k \cup S_t^{ja}$ and $\mathcal{P}_x \setminus \{\mathcal{P}_{x_{itk}}\} : t \in \mathcal{T}$, and $k \in N_{t,f}$, and agent j is re-planned with constraint set $\Omega_{lj} = \Omega_k \cup S_t^{ia}$ and $\mathcal{P}_x \setminus \{\mathcal{P}_{x_{jtk}}\} : t \in N_a$, and $k \in N_{t,f}$. The output trajectory of each agent i from the low-level solver is updated to (ξ_{li}^i) , and pixel coverage set to $\mathcal{P}_{x_{li}}$, such that $P_{li} = (\xi_{li}, \Omega_{li}, \mathcal{P}_{x_{li}}, g_{li})$, where g_{li} is calculated with the updated set of Ω_{li} and $\mathcal{P}_{x_{li}}$. Similarly, for agent “ j ”, new constraints and pixel-set leads to $P_{lj} = (\xi_{lj}, \Omega_{lj}, \mathcal{P}_{x_{lj}}, g_{lj})$. P_{li} , and P_{lj} are added to the TREE with their corresponding g , and the tree is rearranged such that the root node has $\max(g = J(\Xi))$.

Algorithm 2: Single-Agent View Search

Data: ξ_o^i , \mathcal{T} trajectories, envHeightMap(H), agentMaxMotion(m^i), \mathcal{P}_{x_i} , Ω_i

Result: ξ^i

```

1 Initialize  $\xi^i \leftarrow \{\}$ 
2 TREE  $\leftarrow \{\mathbf{R}_0, \xi_o^i, e_0^i, t\}$ 
3 while TREE.top().t is not T do
4    $Q_k \leftarrow (\overline{\mathbf{R}}_k, \xi_k^i, e_{k-1}^i, k)$  TREE.pop()
5   updateProcessedStates[ $\xi_k^i$ ] = true
6   updateStatesFromPrevAction( $Q_k$ )
7    $e_{k_1}^i, \dots, e_{k_{ST}}^i \leftarrow$ 
      availableActions( $\xi_k^i, k, \Omega_i, H, m^i$ )
8   forall  $e_{k_1}^i, \dots, e_{k_{ST}}^i$  do
9      $\overline{\mathbf{R}}_{k+1} \leftarrow (\mathbf{R}_{k+1} + \overline{\mathbf{R}}_k) \gamma^{k-1}$ 
10    // Cumulative reward x discount factor (< 1)
      for encouraging to reach higher reward states
      in lesser timesteps
11     $Q_{k+1} \leftarrow \{(\overline{\mathbf{R}}_{k+1}, \xi_{k+1}^i, e_{k_n}^i, k+1)\}$ 
12    if processedStates[ $\xi_{k+1}^i$ ] == true ||  $\overline{\mathbf{R}}_{k+1} <$ 
        TREE[ $\xi_{k+1}^i$ ].R then
13      continue
14    end
15    TREE  $\leftarrow Q_{k+1}$ 
16  end
17 end
18  $\xi^i \leftarrow \text{BackTracking}(\text{processedStates})$ 
19 return  $\xi^i$ 

```

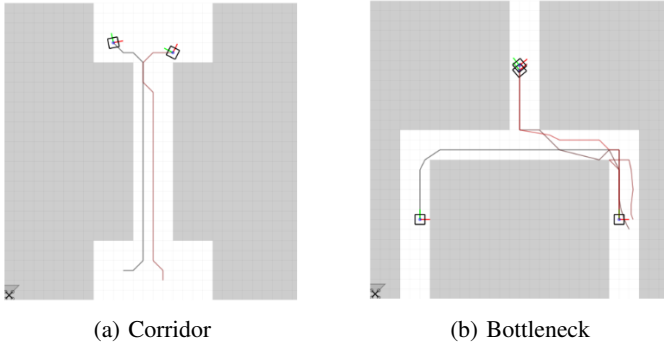


Fig. 4: (a) Corridor and (b) Bottleneck

B. Single-Agent View Planner

The proposed single-agent view search draws from the A^* path planning algorithm, using a perceptual metric as a heuristic to guide node exploration, as shown in Algorithm 2. The algorithm begins at the camera-equipped UAV's start position, adding to $TREE$ all nodes with valid actions within the action space. Each node's reward is adjusted by a discount factor to balance short-term and long-term rewards. A new P_{k+1} is created, and if it's either absent from $TREE$ or has a higher reward than an existing node, it's added to $TREE$. Finally, the trajectory is computed by backtracking from the end position.

VI. EXPERIMENTAL RESULTS

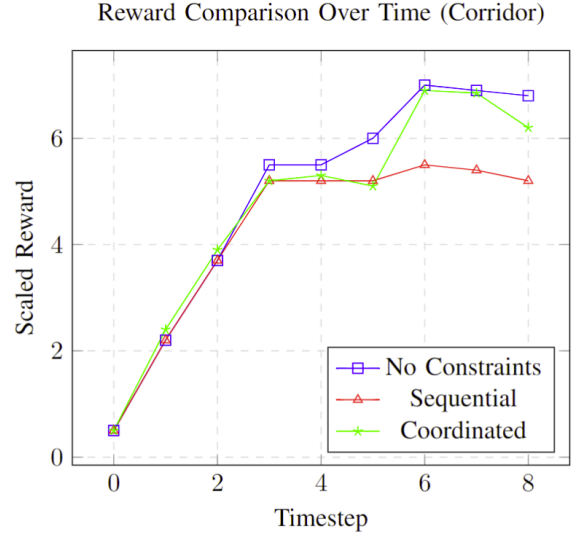
In this section, we present the experimental results of the multi-UAV camera system designed for dynamic actors. The simulations were conducted using a custom-built simulation environment to test the system under bottleneck and corridor scenarios.

A. Ego-Centric Test

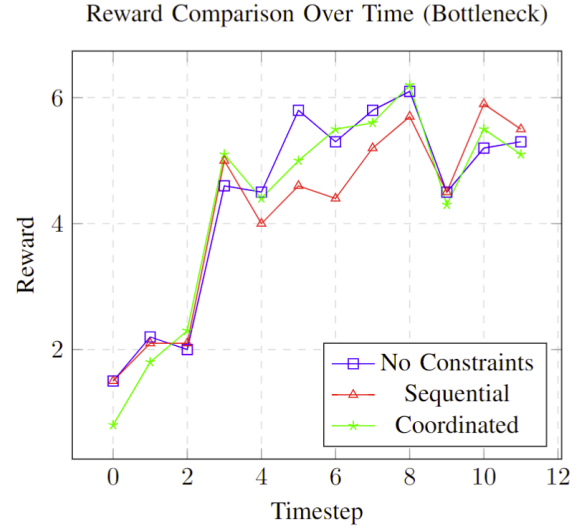
In this test, we evaluated scenarios where the agents exhibited non-egoic behavior, such as one robot getting out of the way to allow another robot to pass as shown in Figure 2. The objective was to assess the system's ability to handle interactions and coordination among agents. We also evaluate the constraints developed in the system when comparing the no-inter-robot constraints, while adding inter-robot constraints, and while doing a conflict-based constraint assignment.

Figure 5a presents the scaled rewards of three view planning methods in a corridor environment with 8 timesteps, 2 robots, and 2 actors. The robots are randomly initialized outside the corridor while the actors move through it in the same direction. From timesteps 3 to 8, when the robots are inside the corridor, a clear gap emerges between the performance of the proposed sequential planning and planning without system constraints. The proposed Alg. 1 aims to narrow this gap, bringing performance closer to the unconstrained system.

Similarly, Figure 5b compares the three methods in a bottleneck environment over 11 timesteps with 4 robots and 4 actors, where actors from different corridors merge in the bottleneck area. From timesteps 3 to 9, the gap between the



(a) Reward comparison amongst different methods in corridor environment.



(b) Reward comparison amongst different methods in bottleneck environment.

Fig. 5: Scale rewards using multiple cameras performing view planning using No inter-robot Constraint, Sequential Constraint, and Conflict Based MDP Value Iteration, over total planning horizon for the environment.

sequential planning view reward and planning without inter-robot constraints highlights the system's reward decrease as constraints are introduced. This indicates the need for an intelligent approach to adding constraints without significantly reducing the view reward. The proposed coordinated capture effectively bridges this gap, achieving performance comparable to the unconstrained method but with added inter-robot constraints.

Together, these two scenarios in Figure 5, with high obstacles and occlusions, demonstrate the proposed methods'

Method	Total Reward (scaled)	Compute Time (s)
Sequential w/ Value Iter.	4127	482
CoCap w/ Value Iter.	4662	3981
CoCap w/ View Search	3922	2.7

TABLE I: Comparison of Total Reward and Compute Time across different methods for corridor.

effectiveness in maintaining high-view rewards despite added constraints.

B. Single-Agent Value Iteration vs. Search

For offline planning approaches, the value iteration solver can compute sub-optimal view positions to achieve high view rewards as shown in GreedyPerspectives [7]. However, in more online settings, there is a need for more efficient algorithms. The proposed single-agent view search, described in Alg. 2, aims to address these use cases.

The Table I compares the total reward (scaled) and compute time (in seconds) for different view planning methods in a corridor environment. The Sequential method with Value Iteration achieved a total reward of 4127 and a compute time of 482 seconds. The Coordinated method with Value Iteration yielded the highest total reward of 4662 but required significantly more compute time, at 3981 seconds. In contrast, the Coordinated method with View Search had a lower total reward of 3922 but was the fastest, with a compute time of only 2.7 seconds. This highlights a trade-off between achieving higher rewards and compute efficiency, showcasing the effectiveness of using single-agent view search in online settings.

VII. CONCLUSION AND FUTURE WORK

In conclusion, this research tackles the challenge of coordinated motion capture of multiple actors in complex, obstacle-dense environments using camera-equipped UAVs. We introduced a novel multi-robot coordinated view planning system, inspired by Conflict-Based Search (CBS) with an occlusion-aware objective. Evaluations in two scenarios demonstrate that our approach, CoCap, outperforms sequential planning, especially in narrow and obstacle-dense environments. Coordinated view planning closely matches the performance of a system without collision constraints, while outperforming sequential greedy planning. Additionally, our single-agent view search method provides a significant computational advantage over the single-agent value iteration solver currently being used.

While promising, the system has limitations. It relies on a 2.5D height map, limiting its handling of overhangs and coverage at doors and windows. Actors are oversimplified as cuboids, and real multi-UAV deployment hasn't been tested, potentially facing communication issues with the centralized planner. More testing is needed in denser, occluded environments, and the search heuristic lacks exploration strategies, especially when high rewards emerge later. Ensuring communication reliability is also crucial for real-world deployment.

ACKNOWLEDGMENT

We would also like to thank Krishna Suresh, Yuechuan Hou, and Micah Nye for their assistance in developing parts of the work. Micah Corah primarily participated in this work while a Postdoctoral Fellow at CMU.

REFERENCES

- [1] J. Gu, T. Su, Q. Wang, X. Du, and M. Guizani, "Multiple moving targets surveillance based on a cooperative network for multi-uav," *IEEE Communications Magazine*, vol. 56, no. 4, pp. 82–89, 2018.
- [2] E. T. Alotaibi, S. S. Alqefari, and A. Koubaa, "Lsar: Multi-uav collaboration for search and rescue missions," *IEEE Access*, vol. 7, pp. 55 817–55 832, 2019.
- [3] S. McCammon, G. Marcon dos Santos, M. Frantz, T. P. Welch, G. Best, R. K. Shearman, J. D. Nash, J. A. Barth, J. A. Adams, and G. A. Hollinger, "Ocean front detection and tracking using a team of heterogeneous marine vehicles," *Journal of Field Robotics*, vol. 38, no. 6, pp. 854–881, 2021. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/rob.22014>
- [4] X. Tong, X. Liu, P. Chen, S. Liu, K. Luan, L. Li, S. Liu, X. Liu, H. Xie, Y. Jin, and Z. Hong, "Integration of uav-based photogrammetry and terrestrial laser scanning for the three-dimensional mapping and monitoring of open-pit mine areas," *Remote Sensing*, vol. 7, no. 6, pp. 6635–6662, 2015. [Online]. Available: <https://www.mdpi.com/2072-4292/7/6/6635>
- [5] Y. Kompis, L. Bartolomei, R. Mascaro, L. Teixeira, and M. Chli, "Informed sampling exploration path planner for 3d reconstruction of large scenes," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 7893–7900, 2021.
- [6] C. Ho, A. Jong, H. Freeman, R. Rao, R. Bonatti, and S. Scherer, "3d human reconstruction in the wild with collaborative aerial cameras," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep. 2021, pp. 5263–5269. [Online]. Available: <https://arxiv.org/pdf/2108.03936>
- [7] K. Suresh, A. Rauniar, M. Corah, and S. Scherer, "Greedy perspectives: Multi-drone view planning for collaborative perception in cluttered environments," Oct. 2024.
- [8] J. Liang and M. C. Lin, "Shape-aware human pose and shape reconstruction using multi-view images," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 4352–4362.
- [9] J. Dong, W. Jiang, Q. Huang, H. Bao, and X. Zhou, "Fast and robust multi-person 3d pose estimation from multiple views," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 7792–7801.
- [10] C. Zheng, W. Wu, C. Chen, T. Yang, S. Zhu, J. Shen, N. Ketharnavaz, and M. Shah, "Deep learning-based human pose estimation: A survey," *ACM Computing Surveys*, vol. 56, no. 1, pp. 1–37, 2023.
- [11] G. Sharon, R. Stern, A. Felner, and N. R. Sturtevant, "Conflict-based search for optimal multi-agent pathfinding," *Artificial Intelligence*, vol. 219, pp. 40–66, 2015. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0004370214001386>
- [12] R. Tallamraju, N. Saini, E. Bonetto, M. Pabst, Y. T. Liu, M. Black, and A. Ahmad, "Aircaprl: Autonomous aerial human motion capture using deep reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6678–6685, October 2020, also accepted and presented in the 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS).
- [13] S. Hughes, R. Martin, M. Corah, and S. Scherer, "Multi-robot planning for filming groups of moving actors leveraging submodularity and pixel density," Milan, Italy, Dec. 2024, to appear.