

Fast Scheduling of Robot Teams Performing Tasks With Temporospatial Constraints

Matthew C. Gombolay, Ronald J. Wilcox, and Julie A. Shah

Abstract—The application of robotics to traditionally manual manufacturing processes requires careful coordination between human and robotic agents in order to support safe and efficient coordinated work. Tasks must be allocated to agents and sequenced according to temporal and spatial constraints. Also, systems must be capable of responding **on-the-fly** to disturbances and people working in close physical proximity to robots. In this paper, we present a centralized algorithm, named “Tercio,” that handles tightly intercoupled temporal and spatial constraints. Our key innovation is a fast, satisficing multi-agent task sequencer **inspired by real-time processor scheduling techniques** and adapted to leverage a **hierarchical problem structure**. We use this sequencer in conjunction with a mixed-integer linear program solver and empirically demonstrate the ability to generate near-optimal schedules for real-world problems an order of magnitude larger than those reported in prior art. Finally, we demonstrate the use of our algorithm in a multirobot hardware testbed.

Index Terms—Human–robot teaming, scheduling.

I. INTRODUCTION

ROBOTIC systems are increasingly entering into domains previously occupied exclusively by humans. In the manufacturing field, for example, there is a strong economic motivation to enable human and robotic agents to cooperatively perform traditionally manual work. This integration requires a choreography of human and robotic work that meets upper and lower bound temporal deadlines for task completion (e.g., the assigned work must be completed within a single shift) and spatial restrictions on agent proximity (e.g., robots must maintain at least 4 m of separation from other agents) in order to support safe and efficient human–robot cooperation. Multi-agent coordination problems with temporospatial constraints can be readily formulated as a mixed-integer linear program (MILP)

Manuscript received February 20, 2017; revised July 11, 2017 and October 19, 2017; accepted October 20, 2017. Date of current version February 5, 2018. This paper was recommended for publication by Associate Editor R. Carloni and Editor T. Murphrey upon evaluation of the reviewers’ comments. This work was supported in part by Boeing Research and Technology and in part by the National Science Foundation Graduate Research Fellowship Program under Grant 2388357. (*Corresponding author: Matthew C. Gombolay.*)

M. C. Gombolay was with Massachusetts Institute of Technology, Cambridge, MA 02139 USA. He is now with Georgia Institute of Technology, Atlanta GA 30332 USA (e-mail: gombolay@csail.mit.edu).

R. J. Wilcox is with Oliver Wyman, New York, NY, USA (e-mail: Ronald.Wilcox@oliverwyman.com).

J. A. Shah is with Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: julie_a_shah@csail.mit.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TRO.2018.2795034

[6], [49]; however, the problem of optimally scheduling $n \geq 3$ tasks (each with a sequence of n_i subtasks) on a set of m machines is NP-hard [73] and is computationally intractable for problems of interest to large-scale factory operations.

Various decentralized or distributed approaches have achieved fast computation and favorable scalability characteristics [8], [14], [16], [68], [76]. Rapid computation is desirable because it enables the system to respond on-the-fly to schedule disturbances, such as an inaccurately performed job, a machine breakdown, or changing temporal constraints [1], [8], [70]. Some approaches boost computational performance by decomposing plan constraints and contributions to the objective function among agents [14]. However, these methods break down when agents’ schedules become tightly intercoupled, such as scenarios wherein multiple agents are maneuvering in close physical proximity to one another. While distributed coordination approaches are necessary for field operations in which environment and geography affect communication among agents, factory operations allow sufficient connectivity and bandwidth for either centralized or distributed approaches to task assignment and scheduling.

In this paper, we present Tercio,¹ a centralized task assignment and scheduling algorithm that scales to multi-agent, factory-size problems and supports on-the-fly scheduling in the presence of temporal and spatial proximity constraints. We empirically demonstrate that this capability enables human and robotic agents to perform manufacturing tasks effectively and in close proximity to one another.

Tercio takes as input a set of tasks composed of precedence-related subtasks; a set of interval temporal constraints relating the start and finish times of subtasks; two-dimensional coordinates specifying the spatial locations where subtasks are performed; physical constraints restricting agent proximity; a set of agent capabilities specifying the tasks and subtasks each agent may perform and the minimum, maximum, and expected time for each agent to complete each task; and an objective function to optimize. Tercio provides a solution consisting of the assignment of agents (i.e., humans or robots) to tasks and a schedule for each agent’s jobs (i.e., start and finish times for each job) such that all temporal and spatial proximity constraints are satisfied and the objective function is empirically within 10% of optimal the majority of the time (see Section X).

¹“Tercio” is named for the Spanish military formation used during the Renaissance period, which consisted of several different types of troops, each with their own strengths, working together as a single unit.

crew assignment and generated Benders cuts for the resulting assignment subproblem. They empirically demonstrated that this method produced more cost-efficient schedules than prior art. Rekik *et al.* similarly employed Benders decomposition to schedule personnel shift-work [66]. The authors applied hand-tailored *forward* and *backward constraints* to cut the search space and proved the correctness of these constraints. Rekik *et al.* showed that Benders decomposition can be used to solve particularly challenging problems in which the *forward* and *backward constraints* do not sufficiently prune the search space.

While Benders decomposition has served as the basis for many state-of-the-art scheduling algorithms, several alternative techniques have also successfully combined MILP and CP approaches. Jain and Grossmann [39] presented an iterative method that first solves a relaxed MILP formulation and then searches for the complete solution using CP. When applied to the scheduling of jobs for machines, the MILP relaxation solves for the assignment of jobs and the CP solves for the schedule. If a solution is identified, the algorithm returns the optimal solution; otherwise, the algorithm infers cuts based on the solution of the MILP relaxation and solves the new MILP relaxation using these cuts. The authors reported results for problems involving up to 20 tasks and 5 machines. Li and Womer later improved on this work by employing a hybrid Benders decomposition algorithm with MILP and CP solver subroutines [50] and reported that their method can solve problems involving 30 tasks and eight different agent types up to four times faster than a standard MILP. Ren and Tang [67] took a similar approach but employed heuristic strategies to generate informative cuts in the event that the CP solver was unable to identify a feasible task sequence. A related method proposed by Harjunkoski and Grossman [37] utilized an iterative approach to producing task assignments and schedules. Both Ren and Tang and Harjunkoski and Grossman empirically demonstrated that their algorithms can solve problems involving up to eight machines and 36 jobs; however, these works did not address problems associated with cross-schedule dependencies.

Although not mathematical programs, some prior works have incorporated planning domain definition language (PDDL)-style problem formulations. For example, Erdem *et al.* [23], [24] developed distributed and semidistributed techniques for scheduling problems involving precedence and absolute temporal constraints, as well as other resource and spatial proximity constraints. These works utilize a formulation that supports causality-based reasoning. The Erdem *et al.* [23], [24] technique yielded optimal solutions for problems with absolute deadlines and complex geometric constraints and readily scaled to problems involving up to eight agents and 80 tasks.

B. Auction and Market-Based Solution Techniques

Auction methods and other market-based approaches to scheduling problems, such as those developed by Ponda *et al.*, Choi *et al.*, and Liu and Shell, also frequently rely upon decomposition of the problem structure [14], [52], [64]. For example, Ponda *et al.* [64] developed a decentralized, market-based solution technique for allocating tasks to agents, given that tasks

are constrained by time windows. Ponda *et al.* [64] employed the consensus-based bundle algorithm (CBBA) [14], in which the objective function and constraints are decomposed by agent so that each agent can quickly solve for the value of its bid on each task. Thus, while the work by [64] represents an added capability for the CBBA, the scalability of the algorithm is strictly worse than the CBBA.

Recently, Nunes and Gini [58] developed the temporal sequential single-item (TeSSI) auction algorithm for decentralized scheduling, which has been shown to outperform the CBBA. TeSSI takes as input a task set in the form of a simple temporal problem (STP); each task has an earliest start and latest finish time (absolute wait and deadline constraints). Constraints relating tasks comprise travel time constraints; resource constraints are not included. Nunes and Gini [58] empirically demonstrated their approach yields improved performance compared with the CBBA [14]. Furthermore, their approach has since been extended to incorporate precedence relations among tasks [54]. However, TeSSI and its variants [54], [58] solve a narrower class of problem than Tercio, in that while they handle some cross-schedule dependencies in the form of precedence relations, they do not handle subtask-to-subtask deadlines (i.e., deadlines constraining the maximum time between the start time of one subtask and the finish time of the second) or resource capacity constraints.

Other techniques solve the task allocation problem efficiently [3], [5], [52], [74], [85], but do not address the sequencing problem. For example, Sung *et al.* addressed multirobot task allocation where each agent maintains a queue of tasks and partial information about other agents' queues [74]. During execution, agents communicate when possible and choose to exchange tasks using a heuristic approach. Sung *et al.* empirically demonstrated that their algorithm can solve problems involving up to six agents and up to 250 tasks; however, the problems did not involve cross-schedule dependencies or task deadlines.

Liu and Shell recently proposed a novel distributed method for task allocation via strategic pricing [52]. Their work builds upon prior approaches to distributed auction algorithms [5], [85], runs in polynomial time, and produces globally optimal solutions. However, this technique does not consider coupling constraints—for example, a problem in which one agent's assignment directly affects the domain of feasible assignments for other agents, as is the case when agents are performing tasks subject to temporal and resource constraints. Chien *et al.* proposed planning methods for a team of robotic rovers to accomplish a set of scientific objectives [13]. The rovers needed to complete a set of tasks where each task required the use of shared, single-access resources (i.e., "shared resources"). The approach by Chien *et al.* [13] uses an iterative-repair centralized planner coupled with an auction algorithm to perform centralized goal allocation and decentralized route planning and goal sequencing. Chien *et al.* benchmarked against a set of randomized problems with three rovers and 12 goals; however, thorough empirical evaluation with an optimal benchmark was not reported.

Lemaire *et al.* approached the problem of allocating UAVs to tasks, represented by a bipartite graph [48]. Here, one set of tasks (UAV navigation) was required to be completed before

the second set (target sensing). The authors first presented a centralized auction solution and reported empirical results for a problem involving 50 tasks and four agents. Next, they described a distributed approach wherein an auctioneer agent assigns the first set of tasks to the multirobot team, and then the second set of tasks is auctioned. This method supports rescheduling in light of dynamic disturbances occurring during task execution; however, the authors did not report empirical results for their distributed method. Sycara *et al.* [75] explored the problem of task allocation and sequencing for a set of agents. In this work, agents were required to share a finite set of resources necessary for task execution. The authors' formulation was distributed in nature and relied upon multiple heuristics to sequentially construct an effective schedule; however, their approach was suboptimal and did not consider deadline constraints relating tasks.

C. Hybrid Solution Techniques

Other hybrid approaches have integrated heuristic schedulers within the MILP solver to achieve better scalability characteristics. For example, Chen *et al.* incorporated depth-first search (DFS) with heuristic scheduling [12]. In this approach, a DFS algorithm sequentially assigns tasks to agents, and a heuristic scheduling algorithm sequences the tasks according to a minimum slack priority. The algorithm also employs heuristics to guide the order in which tasks are assigned to resources during the search. Chen benchmarked their work on problems involving approximately 50 tasks and 10 resources (or agents) using a standard problem database [21], [60].

Alternatively, Castro and Petrovic [9] used a heuristic scheduler to seed a feasible schedule for the MILP with regard to patient procedures conducted within a hospital. This method incorporates a tiered approach to minimize a three-term objective function: First, a heuristic scheduling algorithm generates an initial feasible solution. Next, an MILP is solved for the first term of the objective function using the heuristic solution as a seed schedule. The MILP is solved again to optimize the second objective function term, using the solution from the first MILP as a constraint. This process repeats, but for the third objective function term. The solution time is reduced by sequentially optimizing the objective function terms; however, this approach sacrifices global optimality.

Other approaches perform cooperative scheduling by incorporating Tabu search within an MILP solver [77], or by applying heuristics to abstract the problem to groupings of agents [44]. These hybrid methods are able to solve scheduling problems involving 5 agents (or groups of agents) and 50 tasks in minutes or seconds and address problems incorporating multiple resources, task precedence, and temporal constraints relating task start and end times to the plan epoch time. However, these approaches do not take more general, task-task temporal constraints into consideration.

Cesta *et al.* [10] addressed the problem of project scheduling with time windows by formulating it as a constraint satisfaction problem. In their work, candidate (potentially infeasible) solutions are initially generated using heuristic methods; random and deterministic heuristic techniques are then used to iter-

tively repair any infeasibilities in the problem. Cesta *et al.* [10] noted that randomization is essential to counteract the bias of greedy scheduling heuristics; however, they did not consider wait constraints that create cross-schedule dependencies, nor did they consider shared resources.

D. Metaheuristic Techniques

Successful metaheuristic methods have included simulated annealing (SA) and genetic algorithms (GAs). Davis [18] produced one of the early works applying GAs to job shop scheduling, although many researchers have since followed suit [25], [26], [31], [87]. Recently, Zhang and Wong [88] developed a GA to perform process planning for single-task machines, single-machine tasks, and time-extended scheduling with precedence constraints; however, the formulation did not consider deadline constraints or shared resource constraints.

Researchers have also sought to apply SA techniques to specific scheduling problems [56], [59], [81]. Prior works have combined GAs and SA to further improve upon solution quality [17], [83]. These techniques rely upon a random walk through the space of possible schedules; as the number of steps in the walk (i.e., algorithm iterations) approaches ∞ , the optimal solution will be found.

These solution techniques are typically applied to job shop scheduling problems in which tasks are related through neither tightly intercoupled upper and lower bound temporal constraints nor shared resource constraints [25], [26], [87], [88]. As they are only probabilistically optimal and rely upon random search, metaheuristics can require a large amount of computation time in order to identify and improve schedules when tight upper bound constraints exist.

E. Application to Task Assignment and Scheduling With Temporospatial Constraints

The problem of scheduling—allocating agents to tasks and sequencing those tasks—has been studied in a wide array of works incorporating various solution techniques. However, prior research does not address the need to quickly solve large-scale problems involving tight dependencies among agents' schedules, which can make decomposition problematic. Typically, allowing multiple robots to work closely within the same physical space produces dependencies among the agents' temporal and spatial constraints, leading to uninformative decompositions. We are unaware of prior work that has yielded a solution technique for time-extended scheduling of heterogeneous agents in which each unit of work (i.e., subtask) is related through upper and lower bound temporal constraints without restriction on problem structure (e.g., only some subtasks can be related by certain constraints), where agents must share access to resources (e.g., physical locations) when performing their work.

While we are not aware of prior work focused on our problem definition, we are able to adapt heuristic, metaheuristic, and exact solution techniques to provide an informative empirical validation. In particular, we benchmarked our solution technique, Tercio, against the following techniques:

- 1) We adapted the insertion-heuristic-based TeSSI algorithm [58] to accommodate our class of problems.

$$W_{\langle \tau_i^j, \tau_x^y \rangle} \leq s_x^y - f_i^j \quad \forall W_{\langle \tau_i^j, \tau_x^y \rangle} \in \mathbf{TC} \quad (4)$$

$$D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s} \geq f_x^y - s_i^j \quad \forall D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s} \in \mathbf{TC} \quad (5)$$

$$D_{\tau_i^j}^{\text{abs}} \geq f_i^j \quad \forall D_{\tau_i^j}^{\text{abs}} \in \mathbf{TC} \quad (6)$$

$$f_i^j - s_i^j \geq lb_{\tau_i^j}^a - M(1 - A_{\tau_i^j}^a) \quad \forall \tau_i^j \in \boldsymbol{\tau}, a \in \mathcal{A} \quad (7)$$

$$f_i^j - s_i^j \leq ub_{\tau_i^j}^a + M(1 - A_{\tau_i^j}^a) \quad \forall \tau_i^j \in \boldsymbol{\tau}, a \in \mathcal{A} \quad (8)$$

$$s_x^y - f_i^j \geq M(x_{\langle \tau_i^j, \tau_x^y \rangle} - 1) \quad \forall \langle \tau_x^y, \tau_i^j \rangle \in \boldsymbol{\tau}_R \quad (9)$$

$$s_i^j - f_x^y \geq -Mx_{\langle \tau_i^j, \tau_x^y \rangle} \quad \forall \langle \tau_i^j, \tau_x^y \rangle \in \boldsymbol{\tau}_R \quad (10)$$

$$s_x^y - f_i^j \geq M(1 + x_{\langle \tau_i^j, \tau_x^y \rangle} - A_{\tau_i^j}^a - A_{\tau_x^y}^a) \quad \forall \tau_i^j, \tau_x^y \in \boldsymbol{\tau} \quad (11)$$

$$s_i^j - f_x^y \geq M(2 - x_{\langle \tau_i^j, \tau_x^y \rangle} - A_{\tau_i^j}^a - A_{\tau_x^y}^a) \quad \forall \tau_i^j, \tau_x^y \in \boldsymbol{\tau}. \quad (12)$$

Employing branch-and-bound search to identify the optimal solution in this MILP-based formulation requires $O(2^{|A||\boldsymbol{\tau}|^3})$ just in terms of the integer variables for allocation ($A_{\tau_i^j}^a$) and sequencing ($x_{\langle \tau_i^j, \tau_x^y \rangle}$). Note that the number of possible sequencing permutations is $O(|\boldsymbol{\tau}|!)$, while the number of sequencing variables in this formulation is $O(|\boldsymbol{\tau}|^2)$. Within the manufacturing settings of interest, the number of tasks and subtasks is typically much larger than the number of agents, so the computational bottleneck when solving for a schedule occurs within the sequencing subproblem.

As noted in Section II, other related works have included formulation of similar scheduling problems as MILPs: For example, Korsah *et al.* [43] proposed a general formulation of the instantaneous assignment problem for the XD class. However, our formulation considers time-extended scheduling in which one must also determine how to sequence jobs. Pinto and Grossmann [61] proposed an MILP formulation that bears similarities to ours: both consider assignments of agents to tasks [see (2)], as well as the ordering of tasks for machines (or agents) as given by (11) and (12). However, our MILP formulation includes shared resource constraints [see (9) and (10)], as well as the unique capabilities or production rates of agents [see (7) and (8)]. Other works have considered separate aspects of our problem, such as inclusion of heterogeneous agents (see [78]) or shared resources (see [41]). However, these works do not propose a mathematical formulation that addresses the XD [ST-SR-TA] problem variant with resource constraints.

Tercio approximately solves our problem by producing sub-optimal makespan solutions. Section X demonstrates through empirical evaluation that the produced makespans are within 10% of optimal for the range of problems evaluated. To preserve temporal flexibility at execution, the solution is then reformulated as an STP that is flexibly dispatched [84] (see Section IV-A).



Fig. 2. This figure depicts the Boeing 777 Fuselage Upright Autonomous Build Process. Courtesy: The Boeing Company.

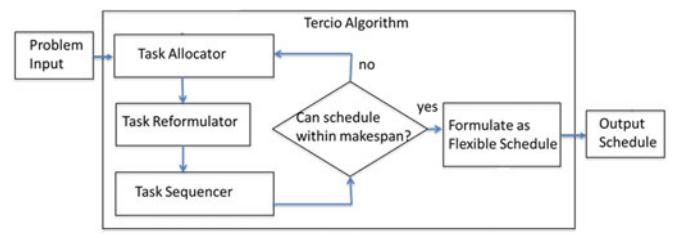


Fig. 3. This figure depicts the system architecture of Tercio.

A. Real-World Motivation Example

Our problem statement is motivated by real-world applications, such as the Boeing 777 Fully Autonomous Upright Build project (see Fig. 2). This type of application requires the coordination of six to eight robots for successful assembly of an aerospace structure. Precedence constraints among the work packages must be respected to ensure structural integrity during the build process. Task–task temporal constraints are imposed by process requirements for the timed application of sealant. The robots’ work must be sequenced to ensure mutually exclusive access to tools, utilities, and floor space; efficient solutions to the problem involve an intricate choreography of robot movements. A fast dynamic scheduling method is necessary to efficiently and effectively respond to schedule disruptions such as those caused by process time variation, rework, and inspection.

IV. TERCIO

In this section, we present Tercio, a centralized task assignment and scheduling algorithm that scales to multi-agent, factory-size problems and supports on-the-fly rescheduling with temporal and spatial proximity constraints. Fig. 3 depicts the system architecture, and the pseudocode for the Tercio algorithm is presented in Fig. 4. Tercio is made efficient through a fast, multi-agent task sequencer inspired by real-time processor scheduling techniques but adapted to leverage the hierarchical problem structure, wherein tasks are composed of precedence-related subtasks. Our approach decomposes the problem statement in Section III into subproblems of task allocation and sequencing. We use an MILP-based solution method (see Section V) to allocate agents to tasks and a polynomial-time task sequencer to efficiently solve the task sequencing problem (see Section VI). Although the task sequencer is satisficing and incomplete, it substantially improves worst-case time complexity.

TERCIO(τ , \mathbf{A} , TC , τ_R , AC , z , cutoff/timeout)

```

1: makespan  $\leftarrow$  inf
2: exclusions  $\leftarrow \emptyset$ 
3: while makespan  $>$  cutoff or runtime  $<$  timeout do
4:    $\mathbf{A} \leftarrow \text{ALLOCATION}(z, \mathbf{A}, \tau, \text{TC}, \text{AC}, \tau_R, \text{exclusions})$ 
5:    $\text{TC} \leftarrow \text{update agent capabilities}$ 
6:   makespan,seq  $\leftarrow \text{SEQUENCER}(\mathbf{A}, \text{TC}, \tau_R)$ 
7:   exclusions  $\leftarrow \text{exclusions} \cup \mathbf{A}$ 
8: end while
9:  $\text{TC} \leftarrow \text{add ordering constraints to enforce seq}$ 
10:  $\text{STP} \leftarrow \text{DISPATCHABLE}(\text{TC})$ 
11: return STP

```

Fig. 4. Pseudocode for the Tercio algorithm.

(We present time complexity analysis for each component of Tercio at the conclusion of their respective description below.)

A. Tercio Pseudocode

Tercio takes the inputs defined in Section III, along with a user-specified makespan *cutoff* intended to terminate the optimization process. This *cutoff* can often be derived from the temporal constraints of the manufacturing process: for example, a user may specify that a provided task set must be completed within an 8-hour shift. Tercio works by iterating through agent allocations until a schedule can be identified that satisfies the maximum allowable makespan for the problem. However, Tercio can also run as an anytime heuristic, terminating once the allotted time has expired.

As shown in Fig. 4, Tercio initializes the makespan (line 1) and previous solution (line 2), and then iterates lines 3–7 to compute a schedule that meets this makespan. A third-party optimizer (Gurobi) solves the agent-allocation MILP (line 4) and returns the agent allocation matrix \mathbf{A} . Interval temporal (TC) constraints are updated according to this matrix by tightening task time intervals (line 5). For example, if a task is originally designated to take between 5 and 15 min, but the assigned robot can complete it in no fewer than 10 min, the interval tightens from [5, 15] to [10, 15].

The agent allocation, the capability-updated TCs, and the spatial map of tasks are then provided as input to the Tercio multi-agent task sequencer (line 6). The task sequencer (described further in Section VI) returns a tight upper bound on the optimal makespan for the given agent allocation, as well as a sequence of tasks for each agent. While this makespan is longer than *cutoff* (or while the algorithm's runtime has not exceeded the specified timeout), the algorithm iterates lines 3–7, each time adding a constraint (line 7) to exclude previously attempted agent allocations. Tercio terminates when either the returned makespan falls beneath *cutoff* or when no solution can be found after iterating through all feasible agent allocations. Note that, for each iteration, a search tree is generated for the agent allocation. We propose preserving the search tree across iterations in future work to reduce computation time, updating it with excluded allocations.

If the cutoff makespan is satisfied, agent and spatial resource sequencing constraints (interval form of $[0, \infty)$) are added to TC (line 9). The resulting STP, composed of the interval temporal

constraints, is compiled into a dispatchable form (line 10) [57], [84], which guarantees that, for any consistent choice of a time point within a flexible window, there exists a solution that can be identified through one-step propagation of interval bounds. The dispatchable form maintains flexibility to increase robustness to disturbances and has been shown to decrease the amount of time spent recomputing solutions in response to disturbances for randomly generated structured problems by up to 75% [84].

V. TERCIO: AGENT ALLOCATION

The Tercio agent allocator performs agent-task allocation by solving an MILP that includes (2), (7), and (8), ensuring that each task is assigned to exactly one agent and that the allocation does not violate the upper and lower bound temporal constraints. In this work, we investigate the objective of minimizing schedule makespan. This corresponds to solving the MILP defined in Section III according to the objective function depicted in (13), where we seek the optimal makespan (i.e., the overall process duration):

$$\min(z), z = g_{\text{MILP}}(\mathbf{A}, \boldsymbol{\tau}) = \max_{\langle \tau_i^j, \tau_x^y \rangle} (f_i^j - s_x^y). \quad (13)$$

However, decomposition of task allocation and sequencing necessitates an objective function that guides the task allocation subroutine toward solutions that are likely to be favorable for the sequencing subroutine. As such, we developed the following objective function [see (14)], comprising three mixed-integer linear terms [see (15) and (16)]. Equation (15) minimizes the maximum work assigned to any one agent, which mitigates situations resulting in a single agent bottle-necking the schedule. Equation (16) minimizes the total amount of work time (i.e., “agent-hours”) by selecting the most efficient agent for each subtask. As such, Tercio’s agent-allocation subroutine maximizes (14) subject to (15) and (16). In our MILP-based task allocation subroutine, binary decision variable $A_{\tau_i^j}^a$ represents the assignment of agent a to subtask τ_i^j . The worst-case time complexity of assigning one of $a = |\mathbf{A}|$ agents to each subtask in $\boldsymbol{\tau}$ via branch-and-bound search is given by $O(2^{|\mathbf{A}||\boldsymbol{\tau}|})$:

$$\min(z), z = \alpha_1 g_1(\mathbf{A}, \boldsymbol{\tau}) + \alpha_2 g_2(\mathbf{A}, \boldsymbol{\tau}) \quad (14)$$

$$g_1(\mathbf{A}, \boldsymbol{\tau}) \geq \sum_{\tau_i^j \in \boldsymbol{\tau}} A_{\tau_i^j}^a \times C_{\tau_i^j}^a \quad \forall a \in \mathbf{A} \quad (15)$$

$$g_2(\mathbf{A}, \boldsymbol{\tau}) \geq \sum_{a \in \mathbf{A}} \sum_{\tau_i^j \in \boldsymbol{\tau}} A_{\tau_i^j}^a \times C_{\tau_i^j}^a. \quad (16)$$

VI. TERCIO: MULTI-AGENT TASK SEQUENCER

The Tercio task sequencer takes the problem defined in Section III as input, along with a set of task assignments provided by the Tercio agent allocator described in Section V. The task sequencer is satisficing, meaning the produced schedule merely satisfies the constraints of the problem [see (2)–(12)] and does not take an objective function as input. The task sequencer returns a valid task sequence if the algorithm can identify one. Tercio schedules tasks in simulation using the dynamic

TERCIO-SEQUENCER(τ)

```

1: reformulate  $D^{s2s}$  and  $W$ 
2:  $t \leftarrow 0$   $\triangleright$  set simulation time to zero
3: ensure temporal feasibility of  $D^{abs}$ 
4: while true do
5:    $availableTasks \leftarrow$  get and sort subtasks ready for execution
      at time  $t$  according to EDF
6:   for all  $\tau_i^j \in availableTasks$  do
7:     if executing  $\tau_i^j$  at  $t$  will not violate temporospatial
      constraints then
8:        $seq \leftarrow$  schedule  $\tau_i^j$ 
9:     end if
10:    end for
11:    $t \leftarrow t + 1$   $\triangleright$  increment simulation time
12:   if all tasks executed then break;
13:   end if
14: end while
15:  $makespan, seq \leftarrow$  extract multi-agent schedule
16: return  $makespan, seq$ 

```

Fig. 5. Pseudocode for the Tercio multi-agent task sequencer.

priority policy earliest-deadline first (EDF), as well as an online schedulability test that guarantees satisfaction of temporospatial constraints for any opportunistic scheduling policy (i.e., a policy that executes a task if one is available to execute). We formulate this schedulability test as a CP problem and determine whether a full, feasible schedule can be developed if subtask τ_i^j is scheduled at time t .

The Tercio task sequencer is inspired by real-time processor scheduling techniques and operates on the structure of a real-time processor scheduling model, called the “self-suspending task model.” The sequencer performs a rapid variant of “edge-checking” (similar to that performed in [45] and [82]), which we call the multiprocessor Russian dolls test. To our knowledge, the Tercio task sequencer is the first real-time scheduling method for multiprocessor systems that tests the schedulability of nonpreemptive, self-suspending tasks in scenarios where multiple tasks have more than one self-suspension and tasks are constrained by shared memory resources.

We now outline the steps taken by the Tercio task sequencer. In Section VII, we discuss the relationship of our problem to prior art in real-time processor scheduling and describe our problem as a real-time processor scheduling problem. We then present how the problem depicted in Section III is reformulated into a real-time processor scheduling problem. Finally, in Section VIII, we present the task sequencer’s priority policy and multiprocessor Russian dolls test.

A. Multi-Agent Task Sequencer Walk-Through

Here, we address the mechanics of our task sequencing algorithm, as depicted in Fig. 5. In line 1, the task set is reformulated (if necessary) into a specific structure that makes the task sequencer operate more efficiently. The task sequencer requires that every subtask τ_i^j involved in an absolute deadline constraint $D_{\tau_i^j}^{abs}$ or subtask-to-subtask deadline constraint $D_{(\tau_a^b, \tau_i^j)}^{s2s}$ have only one predecessor subtask τ_x^y . This means that either $\tau_x^y = \tau_i^{j-1}$, or else there exists a lower bound wait con-

straint $W_{(\tau_x^y, \tau_i^j)}$. In turn, every τ_x^y must have exactly one such predecessor subtask. This recurses until reaching either the epoch subtask τ_0^0 in the case of an absolute deadline constraint, or τ_a^b in the case of a subtask-to-subtask deadline.

In Section VII-C, we prove that the reformulation process preserves correctness—any schedule that satisfies the constraints of the reformulated task set will also satisfy the constraints of the original set. The ability to reformulate depends upon the laxity of the deadlines within the task set and the structure of the constraints. The tighter the deadlines and the more connected the constraint graph, the less able the algorithm will be to reformulate the problem into the structure our schedulability test requires. However, this more restricted structure enables our schedulability test to compute an empirically tight schedulability test in polynomial time.

In line 2, simulation time is initialized to zero. In line 3, the algorithm determines whether the set of absolute deadlines D^{abs} in the reformulated task set is temporally consistent, meaning that the set of agents will be able to successfully schedule against those deadline constraints and their associated spatial constraints. This is determined using the multiprocessor Russian dolls test, as described in Section VIII. After ensuring schedule feasibility due to D^{abs} , the algorithm begins to schedule all subtasks in τ in simulation (lines 4–14). In line 5, the algorithm prioritizes the order in which it attempts to schedule available subtasks according to the EDF dynamic scheduling priority. EDF, commonly used in real-time systems [36], [86], [89] and multi-agent scheduling [11], attempts to execute the task τ_i^j with the earliest (smallest in magnitude) deadline d_i^j first.

Line 6 iterates over all available subtasks τ_i^j and applies the multiprocessor Russian dolls test (line 7) to determine whether a subtask can feasibly be scheduled at time t while satisfying the temporospatial constraints. If feasible, the subtask is scheduled at time t (line 8) and an all-pairs shortest path (APSP) computation [28], [40] is performed to update the temporal constraints. The simulation time is incremented (line 11) and the algorithm continues until all tasks have been scheduled (line 12). Finally, the algorithm returns the makespan and subtask sequence (line 15).

VII. REAL-TIME PROCESSOR SCHEDULING ANALOGY FOR THE TASK SEQUENCER

The design of our informative, polynomial-time task sequencer was inspired by a processor scheduling analogy in which each agent is a computer processor able to perform one task at a time, and a physical location in discretized space is modeled as a shared memory resource accessible by only one processor at a time. Wait constraints (lower bounds on interval temporal constraints) are modeled as “self-suspensions” [46], [69]—times during which a task is blocked while another piece of hardware completes a different, time-durative task.

Typically, assembly manufacturing tasks have more structure (e.g., parallel and sequential subcomponents) and more complex temporal constraints than real-time processor scheduling problems. AI scheduling methods address complex temporal

constraints and gain computational tractability by leveraging a hierarchical structure within the plan [71]. We bridged AI and real-time processor scheduling in order to develop a fast multi-agent task sequencer that satisfies tightly coupled upper and lower bound temporal deadlines and spatial proximity restrictions (shared resource constraints).

The Tercio task sequencer operates on an augmented *self-suspending task model* (defined in Section VII-B) and returns a valid task sequence if the algorithm can identify one. The task sequencer is satisficing and incomplete; however, we have empirically validated that it returns makespans within 10% of the optimal when integrated with the Tercio agent allocation algorithm (see Section X).

Processor scheduling of *self-suspending task systems* has been the focus of much prior work due to the integration of relatively recent hardware and supporting software systems (e.g., GPUs and PPUs) that trigger the external blocking of tasks [20], [46], [69]. Self-susensions can be thought of as lower bound temporal constraints relating tasks: for example, a user might specify that a first coat of paint needs at least 30 min to dry before a second coat may be applied—this 30-minute wait time is a self-suspension of the painting task.

Prior work has computed the uniprocessor schedulability of a task set with single [51], [69] or multiple [32], [33] self-susensions. In our work, we compute the multiprocessor schedulability of a task set in which multiple tasks have more than one suspension and each subtask has a resource constraint. Our approach incorporates a scheduling policy that partially restricts the behavior of the scheduler in order to reduce incidence of multiprocessor schedule anomalies due to self-susensions. Our approach is similar in spirit to prior art that restricted behavior to reduce anomalies that inherently arise from application of uniprocessor scheduling methods to self-suspending task sets [32], [33], [47], [65], [69]. We first introduce our task model in this section. Second, we describe how our task sequencer satisfies temporospatial constraints in Section VIII.

A. Traditional Self-Suspending Task Model

The Tercio task sequencer relies upon a well-formed task model that captures hierarchical and precedence structure within the task network. The basis for our framework is the self-suspending task model, described as

$$\tau_i : \left(W_{\langle \tau_i^0, \tau_i^1 \rangle}, C_i^1, W_{\langle \tau_i^1, \tau_i^2 \rangle}, C_i^2, \dots, C_i^m, T_i, D_i \right). \quad (17)$$

In this model, a set of tasks τ must be processed by the computer. An instance of each task τ_i is *released* (eligible to execute) at every period T_i . The execution of the first subtask of τ_i may be delayed from the epoch start as specified by a lower bound wait constraint $W_{\langle \tau_i^0, \tau_i^1 \rangle}$, called the “phase offset.” For each task, there are m_i subtasks, with $m_i - 1$ self-suspension intervals for each task $\tau_i \in \tau$. We use τ_i^j to denote the j th subtask of τ_i . C_i^j is the expected duration (cost) of τ_i^j . A subtask τ_i^k is released once its predecessor subtask and preceding self-suspension have both finished executing. $W_{\langle \tau_i^{j-1}, \tau_i^j \rangle}$ is the lower bound wait constraint

(or self-suspension) interval relating the end of τ_i^{j-1} and the start of τ_i^j . T_i and D_i are the period and deadline of τ_i , respectively.

B. Augmented Self-Suspending Task Model

The standard self-suspending task model provides a solid basis for describing many real-world processor scheduling problems of interest [20], [46], [69]. Scheduling problems within the manufacturing domain inherently have strong, hierarchical structures that are captured well by the traditional self-suspending task model. However, self-suspending task systems generally assume that processors are homogeneous and do not include more general temporal or resource constraints among tasks and subtasks [51]:

$$\begin{aligned} \tau_i : & \left(W_{\langle \tau_i^0, \tau_i^1 \rangle}, (C_{\tau_i^1}^a, R_i^1), W_{\langle \tau_i^1, \tau_i^2 \rangle}, (C_{\tau_i^2}^a, R_i^2), \dots \right. \\ & \left. (C_{\tau_i^m}^a, R_i^m), T_i, D_i^{s2s}, D_i^{\text{abs}} \right); \left\{ A_{\tau_i^j}^a \right\}, \left\{ W_{\langle \tau_i^j, \tau_x^y \rangle} \right\}. \end{aligned} \quad (18)$$

We first augmented the model to encode the assignment of processors (i.e., agents) to subtasks, denoted by the set of binary decision variables $\{A_{\tau_i^j}^a\}$. The subtask cost $C_{\tau_i^j}^a$ is now dependent upon the capabilities of the agent a processing the subtask. The second augmentation enables a user to specify subtask-to-subtask deadlines $D_{\langle \tau_i^j, \tau_i^k \rangle}^{s2s}$ between the start of τ_i^j and end of τ_i^k [see (19)]; as well as absolute deadlines $D_{\tau_i^j}^{\text{abs}}$ for the absolute finish times of subtasks [see (20)], as shown in (19) and (20). Subtask-to-subtask deadlines are restricted to constrain subtasks within the same task, enabling fast edge-checking in our schedulability test (see Section VIII). In Section VII-C, we describe how general upper bound deadline constraints relating pairs of subtasks are reformulated into this augmented self-suspending task model:

$$D_{\langle \tau_i^j, \tau_i^k \rangle}^{s2s} \geq f_i^k - s_i^j \quad (19)$$

$$D_{\tau_i^j}^{\text{abs}} \geq f_i^j. \quad (20)$$

We define d_i^j as the implicit deadline on the finish time f_i^j of subtask τ_i^j , which is implied by the absolute and subtask deadlines $\mathbf{D}^{\text{abs}} \cup \mathbf{D}^{s2s}$. These deadline constraints provide additional expressiveness to encode binary temporal constraints relating tasks. For instance, these constraints may be used to specify that a sequence of subtasks involving sealant application must be completed within 30 min of opening the sealant container. These types of constraints are commonly included in AI and operations research scheduling models and are vital for modeling many real-world problems [19], [57].

We augmented the model to express subtask-to-subtask wait constraints [see (21)] provided that the subtasks meet certain criteria. To describe this restriction, we first introduce two categories of subtasks: *free* and *embedded*.

Definition 1: A *free subtask* $\tau_i^j \in \tau_{\text{free}}$ is a subtask that does not share a deadline constraint with τ_i^{j-1} . In other words, a subtask τ_i^j is free iff for any deadline $D_{\langle \tau_i^a, \tau_i^b \rangle}$ associated with that

takes a task set as input and returns either a reformulated task set (if the algorithm can find a solution) or null (if no feasible reformulated task set can be identified).

D. Reformulation Pseudocode

Line 1 checks whether any subtask-to-subtask deadlines $D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}$ exist such that there is neither a wait constraint $W_{\langle \tau_i^j, \tau_x^y \rangle}$ nor a set of wait constraints linking τ_i^j and τ_x^y via one or more intermediary subtasks (e.g., $W_{\langle \tau_i^j, \tau_a^b \rangle}$ and $W_{\langle \tau_a^b, \tau_x^y \rangle}$). If there exists such a deadline $D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}$, it is replaced with a new absolute deadline $D_{\tau_x^y}^{abs}$ that implicitly enforces $D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}$ (line 2). Note in line 2, APSP $_{\langle \tau_i^j, \tau_x^y \rangle}$ is the APSP temporal upper bound between τ_i^j and τ_x^y , and that the APSP temporal lower bound between ordered tasks is nonpositive.

For each subtask with its initiation constrained by a deadline and which involves more than one predecessor subtask, the algorithm iterates over the predecessor subtasks (lines 5 and 6). Lines 7 and 8 consider two predecessors at a time τ_a^b and τ_x^y and determine whether the lower bound constraints can be reformulated such that τ_a^b is the predecessor of τ_x^y , or vice-versa (lines 9 and 10). For a subtask τ_i^j with two predecessors τ_a^b and τ_x^y , the algorithm restructures the problem such that τ_a^b is the predecessor of τ_x^y by replacing $W_{\langle \tau_x^y, \tau_i^j \rangle}$ with $W_{\langle \tau_a^b, \tau_x^y \rangle}$, according to

$$W_{\langle \tau_a^b, \tau_x^y \rangle} = \max(W_{\langle \tau_a^b, \tau_i^j \rangle} - C_{\tau_x^y} - W_{\langle \tau_x^y, \tau_i^j \rangle}, W_{\langle \tau_a^b, \tau_x^y \rangle}). \quad (22)$$

In line 11, in the event that both methods of reformulating the problem are feasible, the algorithm employs a “tie-breaker”: If the temporal distance path $\tau_i^j \rightarrow \tau_x^y \rightarrow \tau_a^b$ is shorter than that between $\tau_i^j \rightarrow \tau_a^b \rightarrow \tau_x^y$, the path from τ_i^j to τ_x^y is removed and a temporal wait constraint is added from τ_a^b to τ_x^y —as a result, τ_x^y is now the predecessor of τ_a^b . If only one method of reformulation is feasible (lines 17 and 19), the constraints are modified accordingly (lines 18 and 20). Each time the constraint network is modified, the reformulation process restarts (lines 22 and 23). The algorithm returns null if it is impossible to reorder τ_a^b and τ_x^y and if there are no more predecessor subtasks of τ_i^j to process (line 24); otherwise, the algorithm returns the reformulated task set (line 28).

E. Reformulation Example

In this section, we provide an example to illustrate the reformulation process. Figs. 7 and 8 depict a task set before and after reformulation, respectively. In these figures, subtask start and end times are denoted as nodes (black circles), and constraints are represented by edges. Blue edges indicate lower bound temporal constraints (phase offsets, self-suspensions, or wait constraints), and orange edges represent upper bound temporal constraints (absolute or subtask-to-subtask deadline constraints). Recall that a task set must satisfy two conditions in order to be correctly sequenced by Tercio: First, for every subtask-to-subtask deadline $D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}$,

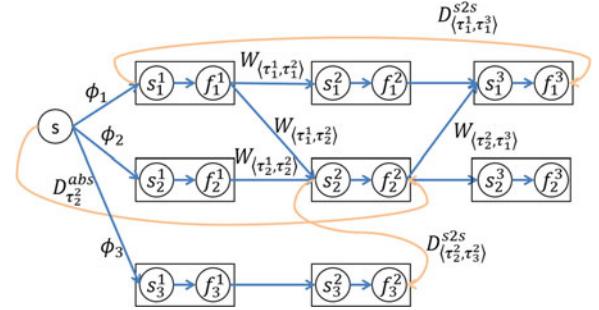


Fig. 7. This figure depicts a task set that must be reformulated to adhere to the augmented self-suspending task model.

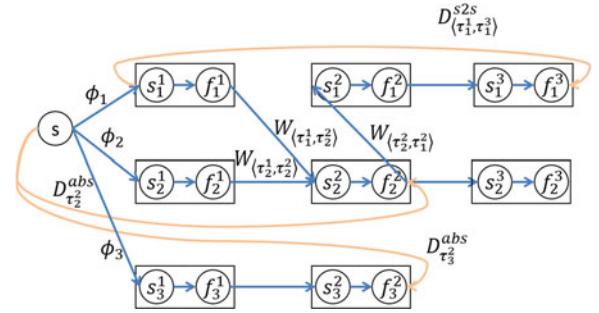


Fig. 8. This figure depicts the first two steps for the reformulation of the task set in Fig. 7.

precedence from τ_i^j to τ_x^y . Second, each subtask must have exactly one predecessor subtask.

Consider the task set depicted in Fig. 7. This example includes three temporal constraints that conflict with the above-mentioned requirements: $D_{\langle \tau_2^y, \tau_3^y \rangle}^{s2s}$ does not satisfy the requirement of a directed path of precedence from τ_2^y to τ_3^y . Also, embedded subtask τ_2^y has more than one precedence constraint delaying its initiation: ($W_{\langle \tau_1^j, \tau_2^y \rangle}$ and $W_{\langle \tau_2^y, \tau_1^j \rangle}$).

The first step to addressing these temporal constraints is to reformulate $D_{\langle \tau_2^y, \tau_3^y \rangle}^{s2s}$ as an absolute deadline $D_{\tau_3^y}^{abs}$, as shown in Fig. 8, such that $D_{\tau_3^y}^{abs}$ is sufficiently tight to guarantee that f_3^2 does not occur later than $s_2^2 + D_{\langle \tau_2^y, \tau_3^y \rangle}^{s2s}$. To ensure consistency, the reformulation subroutine sets $D_{\tau_3^y}^{abs}$ equal to $-APSP(\tau_2^y, \tau_0^0) + D_{\langle \tau_2^y, \tau_3^y \rangle}^{s2s}$. The requirement of a path of precedence from τ_i^j to τ_x^y for every subtask-to-subtask deadline $D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}$ is now satisfied.

Next, we address the wait constraints $W_{\langle \tau_1^j, \tau_2^y \rangle}$ and $W_{\langle \tau_2^y, \tau_1^j \rangle}$. In order to do so, we must first determine whether to require for τ_2^y to precede τ_1^j , or vice-versa. If neither option is possible, the reformulation algorithm cannot reformulate the task set. If only one option is possible, the algorithm selects and enforces that precedence constraint. If both options are temporally feasible, the algorithm assigns the ordering to reduce the delay of τ_2^y . In this example, ordering τ_2^y before τ_1^j is preferred, so the algorithm replaces $W_{\langle \tau_1^j, \tau_2^y \rangle}$ with $W_{\langle \tau_2^y, \tau_1^j \rangle} = \max(W_{\langle \tau_1^j, \tau_2^y \rangle} - (W_{\langle \tau_1^j, \tau_2^y \rangle} C_1^2), 0)$ to ensure satisfaction of the original constraint. The task set can now be correctly sequenced by the Tercio task sequencer algorithm.

Tercio is an iterative algorithm; for this evaluation, we limited the number of iterations to 25. We found that solutions returned by Tercio did not improve significantly with additional iterations for the problems we were able to benchmark against the optimal solution. For agent allocation, we set $\alpha_1 = 2$ and $\alpha_2 = 1$, as this generally achieves helpful allocations for the sequencer, and terminated Gurobi once the incumbent solution was found to be within 0.1% of optimal. We compared our approach against three benchmarks: 1) TeSSI [58], 2) object-coding genetic algorithm (OCGA) [88], and 3) an exact, MILP-based solution method. For all algorithms, a 60-minute timeout was applied.

We compared our approach to a technique based on TeSSI and its variants [54], [58]. At its core, TeSSI is an insertion heuristic. As discussed in [29], insertion heuristics function by deciding which subtask should next be inserted, where, and to which agent based on some prescribed criteria. TeSSI's criterion is the makespan of the agent to which the subtask is assigned. TeSSI operates on a more-restricted problem structure that involves neither upper and lower bound temporal constraints nor resource constraints. Thus, to compute each agent's makespan, TeSSI simply adds the duration of each task and the time spent traveling between tasks, which is linear in the number of subtasks assigned to the agent. In contrast, Tercio considers problems with upper and lower bound temporal constraints and resource constraints. In order to correctly schedule against upper bound temporal constraints, it is necessary to employ a temporal consistency check (e.g., Tercio's Russian dolls test) as commitments are made. A directed path consistency algorithm [19] may also serve this purpose as an alternative to the Russian dolls test. For our comparison, we used Snowball [63], a state-of-the-art algorithm that achieves complexity of $O(n^2)$ for certain cases, such as graphs of constant tree width. We refer to this APSP-variant of TeSSI as TeSSI*.

OCGA is a state-of-the-art GA developed for job shop scheduling problems [88]. We adopted the parameter prescribed by [88] for OCGA as follows: population size $N = 100$, pairs of parents replicating per iteration $r = 7$, probability for crossover $p_c = 1$, probability for shifting genes' loci $p_e = 1$, probability for re-allocating agent for a given subtask $p_m = 0.06$, degeneration ratio $R = 0.09$, and number of iterations $N = 2,500$. As with TeSSI*, it was necessary to add a directed path consistency algorithm to evaluate the quality of each schedule; we also used Snowball for this purpose [63]. In order to improve runtime, we used a hashing function, which stored the quality of previously computed schedules to reduce the number of calls to Snowball with successive iterations of OCGA. GAs such as OCGA are generally probabilistically guaranteed to find the optimal solution as the number of iterations approaches infinity. Thus, given enough time, OCGA will identify a solution better than that generated by Tercio. We find that a more-helpful benchmark than comparing solution quality between OCGA and Tercio is to measure how much time is required for OCGA to find a solution comparable to or better than that identified by Tercio. Thus, in our results, we report the computation time required for OCGA to find such a solution. For instances in which Tercio cannot identify a solution, we report the total

computation time required for all 2500 iterations, as prescribed by [88].

The exact, MILP-based formulation, corresponding to solving the MILP defined in Section III and objective function in (13), was solved by calling Gurobi using its default settings, which sets the optimality threshold to 0.1%.

A. Generating Random Problems

We evaluated the performance of Tercio when applied to randomly generated problems. The lower bound agent-task times were of the form $lb_{\tau_i^j}^a$ and were drawn from a uniform distribution over the interval [1, 10]. Expected agent-task times $C_{\tau_i^j}^a$ were generated from $C_{\tau_i^j}^a \sim U[lb_{\tau_i^j}^a, 10]$. We did not explicitly constrain the maximum amount of time an agent spent performing a subtask. Approximately 25% of the generated precedence constraints $W_{\langle \tau_x^{y-1}, \tau_x^y \rangle}$ were wait constraints with nonzero lower bounds drawn from the interval [1, 10]. Approximately 25% of subtasks had a wait constraint $W_{\langle \tau_i^j, \tau_x^y \rangle}$, in addition to their self-suspension constraint $W_{\langle \tau_x^{y-1}, \tau_x^y \rangle}$.

In order to generate subtask-to-subtask deadlines, we randomly sampled two subtasks τ_i^j and τ_x^y , such that $y \geq j$, and added a deadline constraint $D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}$. To add absolute deadlines, we randomly sampled a subtask τ_i^j and added a deadline constraint $D_{\tau_i^j}^{\text{abs}}$. We sought to generate problems of sufficient challenge for our validation. We established a metric $\hat{D} = \sum_{D_{\langle \tau_i^j, \tau_x^y \rangle}^{s2s}} (y - j + 1) + \sum_{D_{\tau_i^j}^{\text{abs}}} j$ and set it to $\frac{1}{4}$ for our empirical evaluation depicted by Fig. 10(a)–(l). The upper bound of each deadline constraint was drawn from a uniform distribution, with the lower bound set to the tightest feasible deadline and the upper bound set to the sum over subtask costs and wait constraint times. The one-dimensional physical locations of subtasks were drawn from a uniform distribution $[1, |\tau|]$. The number of subtasks m_i within each task τ_i was drawn from a uniform distribution in the interval $[1, 2n]$, where n is the number of tasks τ_i in τ . While an agent performed subtask τ_i^j at resource location (x) , no other agent could work on a subtask at location (u) if $x - 1 \leq u \leq x + 1$.

B. Computation Speed and Scalability

Fig. 10(a), (e), and (i) depict our evaluation of Tercio's scalability and computational speed. We present the median and upper/lower quartiles of the computation time for 50 randomly generated problems incorporating 5, 10, or 100 agents [see Fig. 10(a), (e), and (i), respectively] and between 4 and 1024 subtasks. Where possible, we provide the computation time for the MILP-based solution method, the insertion algorithm, and the time required for OCGA to find a solution as good as that identified by Tercio. Note that the median is typically reported for such optimization problems because the distributions are often skewed and the mean is a less informative measure [79].

Tercio can solve problems involving up to 100 agents and 1000 subtasks in ~ 120 s, a substantial improvement over our benchmarks. We note that TeSSI* and OCGA compute more

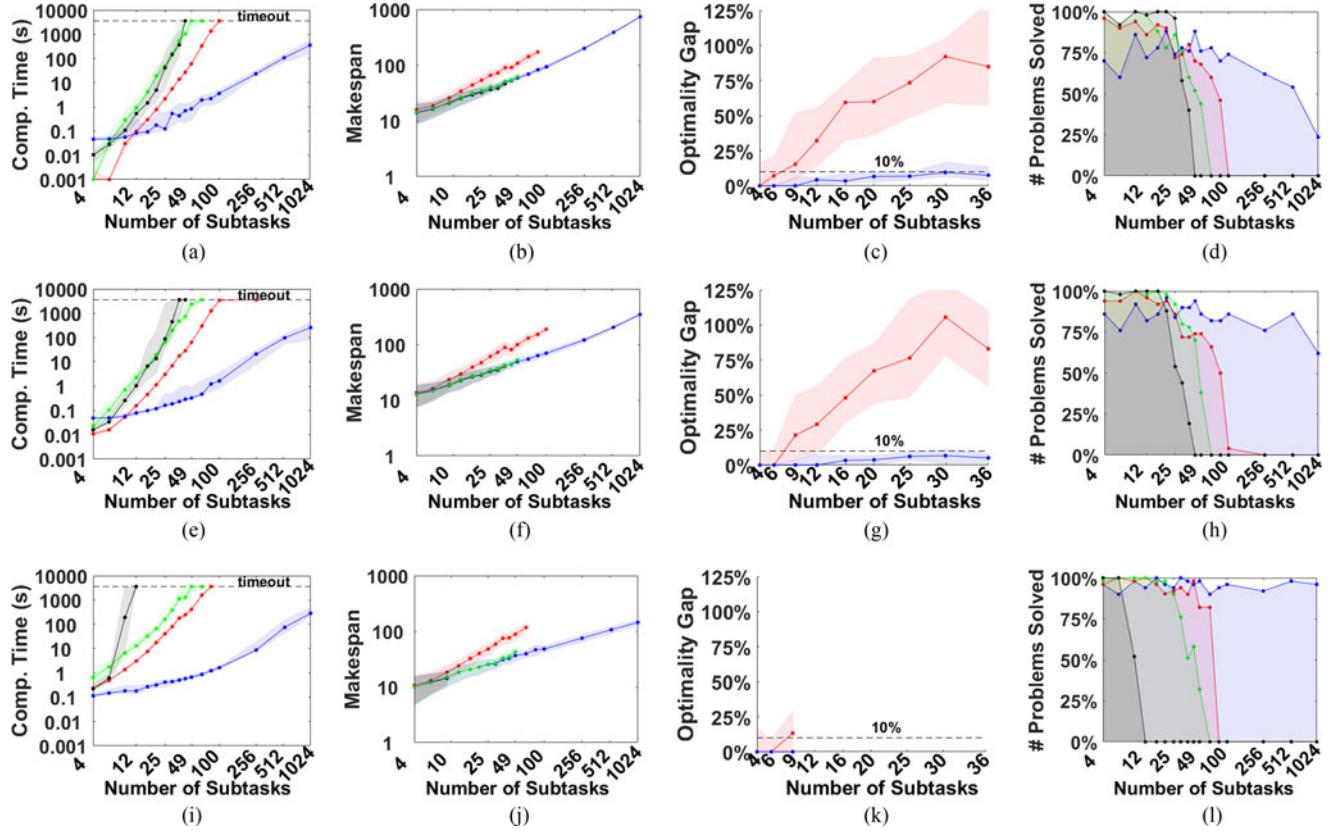


Fig. 10. This figure depicts the performance of Tercio relative to TeSSI*, OCGA, and the exact solution (where possible). (a)–(d), (e)–(h), and (i)–(l) present results for problems involving 5, 10, and 100 agents, respectively. When reporting computation time or solution quality (i.e., the makespan), the graphs depict the median and quartiles. Tercio is represented in blue, TeSSI* in red, OCGA in green, and the MILP-based solution technique in black. Note that OCGA is not depicted when reporting the optimality gap, as we terminated OCGA once it found a solution as good as that of Tercio.

slowly than previously published variants because of the need to include a temporospatial feasibility test for each assignment according to the upper and lower bound temporospatial constraints, requiring $O(n^2)$ step in the innermost loop of the algorithm. This feasibility test greatly increases computation time despite using the fastest technique available [63]. Our findings underscore the benefit to computation time provided by our schedulability test.

C. Optimality

We empirically validate that Tercio produces solutions within 10% of optimal. Fig. 10(c), (g), and (k) depict the median and upper/lower quartiles of the makespans produced by Tercio, along with the insertion algorithm for 50 randomly generated problems involving 5 and 10 agents and problem sizes spanning four to 42 tasks. The median deviation from optimal for Tercio was less than 10% for all testable problem sizes. However, the median deviation from optimal for the insertion algorithm increased up to 100% for larger problems.

D. Evaluating Completeness

Fig. 10(d), (h), and (l) depict the proportion of problems solved by Tercio with 5, 10, and 100 agents, respectively. While TeSSI* and OCGA were able to solve slightly more smaller problems than Tercio, Tercio's ability overtook that of both

TeSSI* and OCGA as problem size increased. Furthermore, Tercio's completeness is proportional to the number of agents; conversely, OCGA and TeSSI* appear less able to identify satisfactory solutions as the number of agents increases.

E. Robustness

To test the robustness of our approach, we considered problems that are more- or less-constrained in Fig. 11(a)–(d), (e)–(h), and (i)–(l). We set $\hat{D} = \{\frac{0}{4}, \frac{1}{4}, \frac{2}{4}\}$. Furthermore, we varied the proportion of subtasks with wait constraints $W_{\langle \tau_i^j, \tau_x^y \rangle}$ such that $\tau_x^y \neq \tau_i^{j+1}$ in the set $\{\frac{0}{4}, \frac{1}{4}, \frac{2}{4}\}$. We found that, relative to OCGA and TeSSI*, Tercio's performance remained strong across a range of constraint settings. Tercio's completeness did degrade for the most-constrained problems [see Fig. 11(l)], but its scalability and solution quality remained strong relative to our benchmarks.

F. Robot Demonstration

Here, we demonstrate the use of Tercio in two hypothetical manufacturing scenarios. In both cases, a team (i.e., set) of robots worked to complete tasks on a simulated fuselage. The robots performed their tasks at specific locations on the factory floor, where there can be multiple subtasks at each location. In order to prevent collisions, each robot reserved the physical location for its subtasks, along with any immediately adjacent

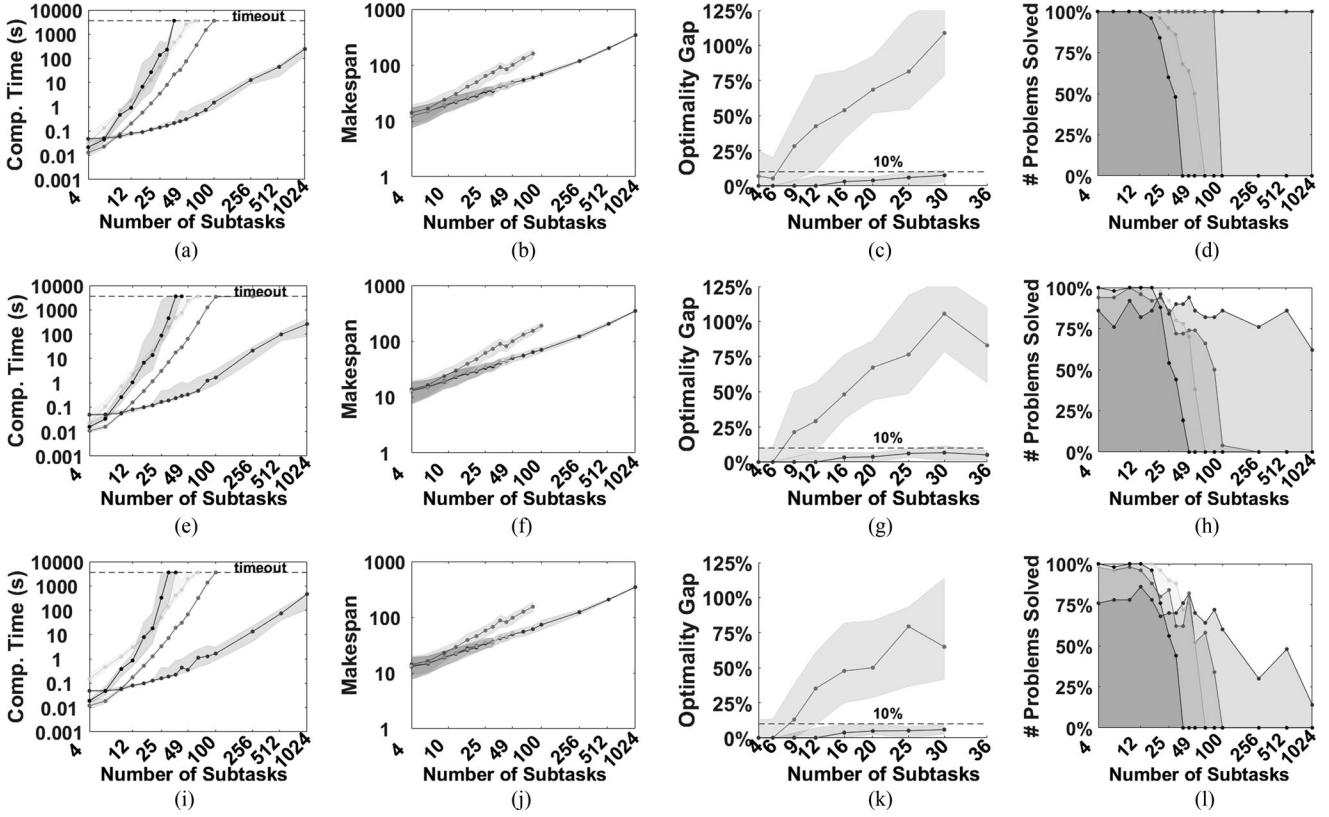


Fig. 11. This figure depicts the performance of Tercio relative to the exact solution, TeSSI*, and OCGA. (a)–(d), (e)–(h), and (i)–(l) present results for problems that involve ten agents and are lightly, moderately, and severely constrained, respectively. Tercio is depicted in blue, TeSSI* in red, OCGA in green, and the MILP-based solution technique in black.

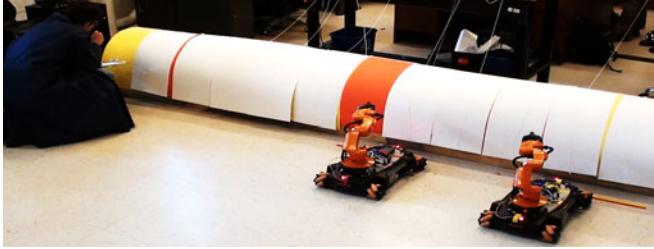


Fig. 12. Robotic team completes mock aerospace final-assembly tasks while maintaining a safe distance from human workers.

task locations; thus, no two workers could be present at the same location (or in neighboring locations) at the same time. For simplicity, only absolute deadlines D^{abs} were considered, although other constraints could be easily incorporated.

In the first evaluation, two KUKA Youbots simulated completion of drilling tasks on an aerospace fuselage, as shown in Fig. 12 (video available at <http://tiny.cc/t6wjxw>). Initially, the robots planned to evenly split 12 identical tasks down the middle of the fuselage. After the robots completed their first subtasks, a worker then requested time to inspect the completed work along the left half of the fuselage; in the problem formulation, this corresponds to adding a resource reservation for the left half for a specified period of time. Tercio replanned in response to the addition of this new constraint and reallocated work in a reasonable manner to make productive use of both

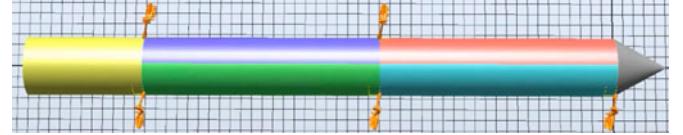


Fig. 13. Robotic team adapts performance of work related to a virtual fuselage in response to multiple disruptions.

robots. For the purposes of this demonstration, the replanning process involved calling Tercio while forcing a fixed assignment and ordering of all subtasks completed prior to the replan request. More-efficient methods could also be considered, such as pruning completed subtasks and reformulating appropriate temporal constraints.

Second, we demonstrated Tercio on a larger simulated problem, as shown in Fig. 13 (video available at <http://tiny.cc/jladjy>). In this demonstration, which incorporated ABB's Robot Studio simulation environment, Tercio coordinated 5 robots to perform 110 identical tasks around an aerospace fuselage. Tercio was applied to replan in response to three disturbances: 1) a human worker's request to enter the space to perform a quality-assurance (QA) inspection, 2) a robot breakdown, and 3) changing task deadlines. The duration of the QA request and the robot breakdown were known at the time of the disturbance. Tercio modeled the QA request as a resource reservation for the section of the fuselage to be inspected (i.e., no robots could occupy that

- [82] P. Vilím, R. Barták, and O. Čepek, “Extension of $o(n \log n)$ filtering algorithms for the unary resource constraint to optional activities,” *Constraints*, vol. 10, no. 4, pp. 403–425, 2005.
- [83] L. Wang and D.-Z. Zheng, “A modified genetic algorithm for job shop scheduling,” *Int. J. Adv. Manuf. Technol.*, vol. 20, no. 1, pp. 72–76, 2002.
- [84] R. J. Wilcox, S. Nikolaidis, and J. A. Shah, “Optimization of temporal dynamics for adaptive human-robot interaction in assembly manufacturing,” in *Proc. Robot. Sci. Syst.*, 2012, pp. 441–448.
- [85] M. M. Zavlanos, L. Spesivtsev, and G. J. Pappas, “A distributed auction algorithm for the assignment problem,” in *Proc. 2008 47th IEEE Conf. Decis. Control*, 2008, pp. 1212–1217.
- [86] F. Zhang and A. Burns, “Schedulability analysis for real-time systems with EDF scheduling,” *IEEE Trans. Comput.*, vol. 58, no. 9, pp. 1250–1258, Sep. 2009.
- [87] F. Zhang, Y. Zhang, and A. Y. C. Nee, “Using genetic algorithms in process planning for job shop machining,” *IEEE Trans. Evol. Comput.*, vol. 1, no. 4, pp. 278–289, Nov. 1997.
- [88] L. Zhang and T. Wong, “An object-coding genetic algorithm for integrated process planning and scheduling,” *Eur. J. Oper. Res.*, vol. 244, no. 2, pp. 434–444, 2015.
- [89] H. Zhao, L. George, and S. Midonet, “Worst case response time analysis of sporadic task graphs with EDF non-preemptive scheduling on a uniprocessor,” in *Proc. 3rd Int. Conf. Autonomic Auton. Syst.*, 2007, pp. 22–22.



Matthew C. Gombolay received the B.S. degree in mechanical engineering from the Mechanical Engineering Department, Johns Hopkins University, Baltimore, MD, USA, the S.M. degree in aeronautics and astronautics from the Department of Aeronautics and Astronautics, Massachusetts Institute of Technology (MIT), Cambridge, MA, USA, and the Ph.D. degree in autonomous systems from MIT, in 2011, 2013, and 2017, respectively. He will join the faculty in the School of Interactive Computing, Georgia Institute of Technology, Atlanta, GA, USA, as the Catherine M. and James E. Allchin Early-Career Associate Professor in Fall 2018, where he will develop novel computational methods for dynamic, fluent human–robot teaming.



Ronald J. Wilcox received the B.A. degree in physics from College of William and Mary, Williamsburg, VA, USA, and the M.S. degree in mechanical engineering from the Mechanical Engineering Department, Massachusetts Institute of Technology, Cambridge, MA, USA, in 2011 and 2013, respectively.

He is currently an Associate with Oliver Wyman, New York, NY, USA



Julie A. Shah received the S.B. and S.M. degrees in aeronautics and astronautics from the Department of Aeronautics and Astronautics, Massachusetts Institute of Technology (MIT), Cambridge, MA, USA, and the Ph.D. degree in autonomous systems from MIT, in 2004, 2006, and 2010, respectively.

She is an Associate Professor of aeronautics and astronautics with MIT and the Director of the Interactive Robotics Group, which aims to imagine the future of work by designing collaborative robot teammates that enhance human capability. Before joining the faculty, she worked with Boeing Research and Technology on robotics applications for aerospace manufacturing.