

MOTIVATING EXAMPLE

- **Nereus problem:** Naval anti-air warfare problem, depicted below.
- **Objectif:** Find a resource assignment strategy, that avoids both:
 1. The ship being hit by incoming enemy missiles;
 2. The bounded and shared resources being over-utilized.
- **Main difficulties posed by Nereus:**
 1. An arbitrary number of enemy missiles may occur over periods;
 2. The Lack of the joint-model of rewards and transitions of the team of enemy missiles;
 3. The large number of resources to allocate to the team of enemy missiles;
 4. The availability of resources positive and negative interactions;
 5. The resource constraints including tight response-time constraints.



FIGURE 1: Naval Environment for Resource Engagement in Unpredictable Situations.

CYCLIC PROGRESSIVE REASONING UNIT (C-PRU)

- **Model of Task Structure:**
 1. Each enemy missile is formalized as a structured task;
 2. Each task describes the ship engagement procedure as illustrated below.

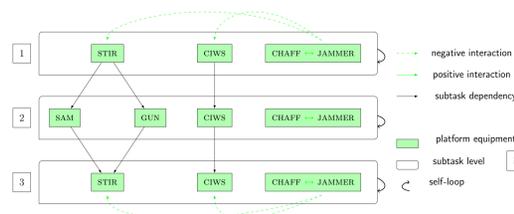


FIGURE 2: Cyclic Progressive Reasoning Unit of a missile engagement procedure.

- **More Formally,** a C-PRU $C_i: \langle S_i, A_i, T_i, R_i, K_i, U, \varphi, \lambda \rangle$ consists of:
 1. A single agent Markov decision problem: $\langle S_i, A_i, T_i, R_i, \lambda \rangle$;
 2. A resource management problem: $\langle K_i, U, \varphi \rangle$:
 - K_i is the set of resource types, such as CHAFF CLOUD or JAMMER as depicted above.
 - $\varphi(a_i, \Omega_i) \in [0, 1]$ is the discount factor of executing action a_i , when actions Ω_i are operating. Negative and positive interactions are illustrating in Figure above.
 - $U = [u_k]_k$ available amount of resources per type $k \in K_i$.
- **Remark:** C-PRU extends the progressive reasoning unit model first introduced by [Mouaddib & al.], for solving infinite-horizon MDPs handling structured tasks.

TWO PHASES APPROXIMATE STRATEGY

- **Inconvenient:** A periodic Multi-agent C-PRUs problem is PSPACE-HARD.
- **Solution:** Computing an approximate solution that sacrifices optimality for computational feasibility.

1. **Offline phase:** Compute heuristic estimates of individual value function of each task:

$$q_i(s_i a_i | \Omega_i) = R_i(s_i a_i | \Omega_i) + \lambda \sum_{s'_i} T_i(s_i a_i s'_i) V_i^*(s'_i | \emptyset) \quad (1)$$

$$v_i(s_i | \Omega_i) = \max_{a_i} q_i(s_i a_i | \Omega_i)$$

where $V^*(s'_i | \emptyset)$ is an optimistic value when ignoring resource interactions, and

$$R_i(s_i a_i | \Omega_i) = \sum_{s'_i} \varphi(a_i, \Omega_i) T_i(s_i a_i s'_i) R_i(s_i a_i s'_i) \quad (2)$$

2. **Online phase:** Recover an approximate estimate of the exact value function of the team of enemy missile tasks.

$$\bar{Q}(sa | \Omega) = \sum_{k=1}^{|M|} \zeta_k \prod_{i \leq k} q_i(s_i a_i | \Omega_i) \quad (3)$$

where ζ_k stands for the weight associated with task M_k , M is the set of all task M_k . The objective consists in finding a joint-action \bar{a}^* at each time period such that:

$$\bar{a}^* = \arg \max_{a \in A} \bar{Q}(sa | \Omega) \quad (4)$$

- **Challenges:** To find a near-optimal solution through a search space bounded by $O(|\otimes_i A_i|^M)$ with respect to the response-time constraints: Each iteration improves the current complete solution.

REAL-TIME DYNAMIC A* (RTDA*)

- **Search Space:** Tree search structure $Tree = (N, E)$
- $N = \{\sigma_i^j\}_{i,j}$ denotes the set of nodes $\sigma_i^j = \langle s_i, \Omega_i^j, Q_i^j \rangle$ describing agent M_i characteristics, where Q_i^j states the heuristic estimate of the current best joint action $\langle a_i \dots a_{|M|} \rangle$:

$$Q_i^j = \sum_{i'=i}^{|M|} \zeta_{i'} \prod_{k \leq i' \leq j} q_k(s_k a_k | \Omega_k^j) = q_i(s_i a_i | \Omega_i^j) (\zeta_i + Q_{i+1}^j) \quad (5)$$

- **RTDA* Algorithm:**

Require: agent states $\{s^j\}_{j=1}^{|M|}$ ordered following EDF.

Ensure: *tree*.

- 1: *tree* \leftarrow EmptyTree
- 2: *open* \leftarrow EmptyStack
- 3: SEARCHSUCCESSORS(σ_1^0).
- 4: **while** *open* \neq EmptyStack **do**
- 5: $\langle \sigma_i^j, a_i^j, \sigma_{i+1}^j \rangle \leftarrow$ *open.pop()*.
- 6: **if** σ_{i+1}^j is not yet *visited* **then**
- 7: add $\langle \sigma_i^j, a_i^j, \sigma_{i+1}^j \rangle$ to *tree*.
- 8: SEARCHSUCCESSORS(σ_{i+1}^j).
- 9: **else**
- 10: update Q_i^j as mentioned eq. (5).
- 11: **end if**
- 12: **end while**

3 AGENTS TOY EXAMPLE

- **Branch-And-Bound Algorithm:** RTDA* uses an upper U_b and lower L_b bounds to prune dominated nodes:

$$U_b(a_i^j) = q_i(\zeta_i + v_{i+1}(\zeta_{i+1} + \zeta(|M| - i - 1))) \quad (6)$$

$$L_b(\sigma_i^j) = Q_i^j$$

where $q_i = q_i(s_i a_i | \Omega_i^j)$, $v_i = v_i(s_i | \Omega_i^j)$ and $\zeta = \max_i \zeta_i$. As an example, action a_2^0 estimate value for $\zeta_i = \frac{1}{3} (\forall i)$, is given by: $L_b = Q_2^0 = .34 \times (\frac{1}{3} + .44) = .2629$. The *upper bound* of selecting action a_2^0 (resp. a_2^1), is given by:

$$U_b(a_2^0) = .23 \times (\frac{1}{3} + .44 \times (\frac{1}{3} + \frac{1}{3}(3 - 2 - 1))) = .1104$$

Thus it is not necessary to expand either node σ_3^0 or σ_3^2 because of $U_b(a_2^0) < L_b$.

- **Anytime Algorithm:** RTDA* selects in a greedy fashion its joint-actions based on pre-compiled solutions.

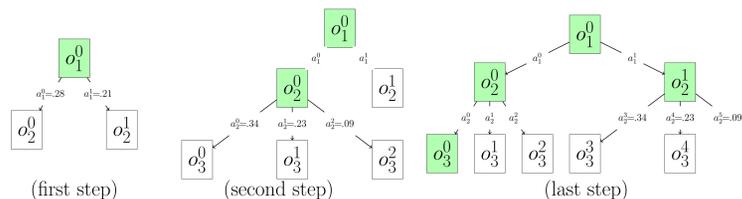
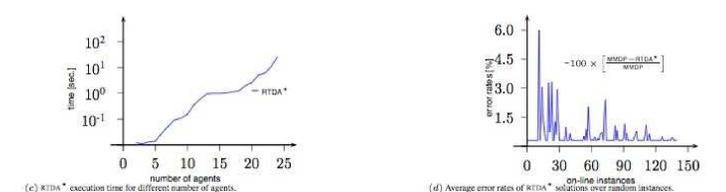
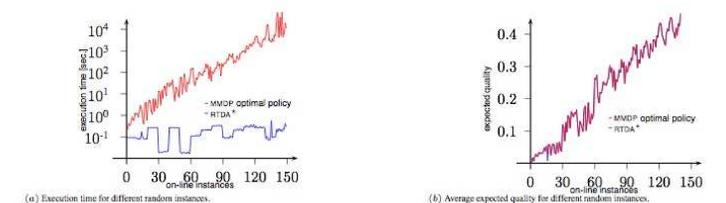


FIGURE 3: Each edge is labeled with an admissible action and its pre-compiled value. The light-green nodes mark the partially or completely explored node.

EXPERIMENTAL RESULTS AND CONCLUSION

- **Experiments:** Comparison of RTDA* and optimal MMDP solutions.



- **Results:** RTDA* is able to provide near-optimal solutions, with respect to bounded and shared resource, and under a large number of agents.

- **Future work directions:**

1. Extend RTDA* in order to find optimal solutions;
2. Adapt RTDA* for solving Dec-(PO)MDPs handling a large number of agents.